

**O'REILLY®**  
Technical Guide

# The LLM Mesh

An Architecture for Building Agentic  
Applications in the Enterprise

Compliments of



**dataiku**

**Kurt Muehmel**



# Create, Connect, and Control AI Agents at Scale With Dataiku

With Dataiku, build analytics, machine learning models, *and* agents all within the same platform so that every new data insight or model created can directly enhance your AI agent efforts.

## ➤ BUILD AI AGENTS YOUR WAY

Design and develop enterprise-grade AI agents with full-code flexibility or a code-free visual interface.

## ➤ STREAMLINE TOOL MANAGEMENT

Manage all your AI agent tools in one place, so you can focus on innovation instead of integration.

## ➤ TROUBLESHOOT AI AGENTS WITH CONFIDENCE

Quickly debug agent behavior with a visual tool that brings full transparency to agent actions.

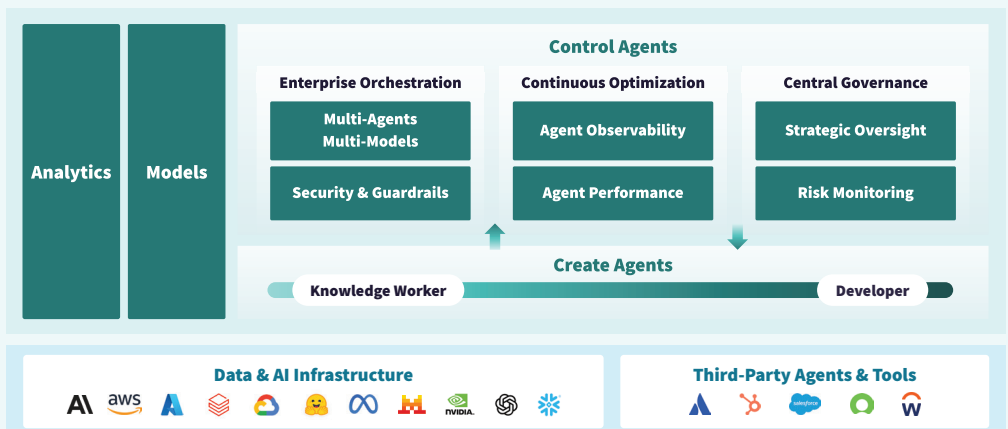
## ➤ DEPLOY & ORCHESTRATE WITH EASE

Put AI agents to work faster and at scale with centralized deployment, smart routing, and built-in performance oversight.

## ➤ UNLEASH AGENTS WITH THE RIGHT GUARDRAILS & GOVERNANCE

Keep agents safe and compliant with integrated LLM security, usage controls, and enterprise-grade governance at every step.

### Dataiku Orchestrates Analytics, Models, & Agents



MORE ON AI AGENTS WITH DATAIKU

---

# The LLM Mesh

*An Architecture for Building Agentic  
Applications in the Enterprise*

*Kurt Muehmel*

O'REILLY®

## The LLM Mesh

by Kurt Muehmel

Copyright © 2026 O'Reilly Media, Inc. All rights reserved.

Published by O'Reilly Media, Inc., 141 Stony Circle, Suite 195, Santa Rosa, CA 95401.

O'Reilly books may be purchased for educational, business, or sales promotional use. Online editions are also available for most titles (<http://oreilly.com>). For more information, contact our corporate/institutional sales department: 800-998-9938 or [corporate@oreilly.com](mailto:corporate@oreilly.com).

**Acquisitions Editor:** Aaron Black  
**Development Editor:** Jeff Bleiel  
**Production Editor:** Kristen Brown  
**Copyeditor:** nSight, Inc.

**Proofreader:** Kim Cofer  
**Interior Designer:** David Futato  
**Cover Designer:** Susan Brown  
**Illustrator:** Kate Dullea

December 2025: First Edition

### Revision History for the First Edition

2025-12-15: First Release

The O'Reilly logo is a registered trademark of O'Reilly Media, Inc. *The LLM Mesh*, the cover image, and related trade dress are trademarks of O'Reilly Media, Inc.

The views expressed in this work are those of the author and do not represent the publisher's views. While the publisher and the author have used good faith efforts to ensure that the information and instructions contained in this work are accurate, the publisher and the author disclaim all responsibility for errors or omissions, including without limitation responsibility for damages resulting from the use of or reliance on this work. Use of the information and instructions contained in this work is at your own risk. If any code samples or other technology this work contains or describes is subject to open source licenses or the intellectual property rights of others, it is your responsibility to ensure that your use thereof complies with such licenses and/or rights.

This work is part of a collaboration between O'Reilly and Dataiku. See our [statement of editorial independence](#).

978-1-098-17658-7

[LSI]

---

# Table of Contents

<b>1. Using LLMs in the Enterprise.....</b>	<b>1</b>
What Is an LLM Mesh?	3
The Right Model for the Right Application	6
Bottom Line: Why the LLM Mesh?	20
<b>2. Objects for Building Agentic Applications.....</b>	<b>23</b>
The Potential of New Agentic Applications	24
LLM Mesh-Related Objects: An Overview	29
The Objects of an LLM Mesh in Detail	32
Cataloging LLM-Related Objects	48
Conclusion	50
<b>3. Quantifying and Optimizing the Cost of LLMs in the Enterprise..</b>	<b>51</b>
Quantifying the Costs of Agentic Applications	52
Techniques for Limiting Costs	66
Cost-Efficient AI Operations in the Enterprise	72
Conclusion	76
<b>4. Measuring and Monitoring the Performance of     Agentic Applications.....</b>	<b>77</b>
Measuring and Monitoring the Quality of Generated Text	81
Measuring and Monitoring the Speed of Agentic Applications	98
Conclusion	100

<b>5. Safety.....</b>	<b>101</b>
Unsafe Behavior in Agentic Applications	103
Layered and Reusable Safety Filters	105
Hallucination Detection	109
Explainability and Transparency	111
Adversarial and Robustness Testing	114
Human Feedback as the Ultimate Arbiter of Safety	116
Conclusion	118
<b>6. Agentic Application Security.....</b>	<b>119</b>
A Paradigm Shift: Security in an LLM Mesh	120
Access Control Frameworks	121
Secure Gateway and API Architecture	124
Permissions Management	126
Audit Logging and Traceability	127
Anomaly Detection and Monitoring	129
Secure Deployment Methods	130
Aligning Agentic Application Security with Enterprise Standards	133
Conclusion	134
<b>7. An LLM Mesh in Action: An Example Agentic Application.....</b>	<b>135</b>
The Business Problem: Taming Complexity in Financial Crime Detection	136
Building Within an LLM Mesh Architecture: A Practical Walkthrough	138
The Benefits of an LLM Mesh for Building, Deploying, and Maintaining the AML Assistant	143
Conclusion: From Ad Hoc AI to a Strategic Capability	144

# Using LLMs in the Enterprise

*May you live in times of rapid technological progress.* This is the blessing and the curse of our current moment. Recent advances in AI, plus the release of wildly popular consumer products, have led to a frenzy of interest in, and use of, AI, and large language models (LLMs) in particular, in the enterprise.

However, AI and LLMs remain nascent in the enterprise, meaning that best practices for their use are still being defined. At the same time, the core technologies—the models themselves, technologies to host and serve the models, etc.—are evolving rapidly.

**Table 1-1** provides a brief timeline of the release of various models and technologies that could be relevant for enterprise use.<sup>1</sup> The diversity and speed of release create both opportunities and challenges when you are looking to use these technologies in production use cases.

---

<sup>1</sup> For additional model comparisons, including performance on benchmarks, pricing, and other dimensions, **Vellum's LLM Leaderboard** is a popular reference.

---

*Table 1-1. A (nonexhaustive) timeline of enterprise-relevant model and product releases*

Developer or provider	Model or product	Release date	Description
OpenAI	GPT-3	May 2020	175-billion-parameter LLM with 2,048 token context window
OpenAI	ChatGPT	November 2022	Consumer chatbot application, powered by GPT-3.5
Microsoft Azure	OpenAI Service	January 2023 (public preview)	Managed service offering LLMs from OpenAI
Amazon Web Services (AWS)	Amazon Bedrock	September 2023	Managed service offering LLMs from various developers
Dataiku	LLM Mesh	September 2023	Commercial LLM Mesh offering for connecting to LLMs and building agentic applications in the enterprise
Databricks	DBRX	March 2024	Open-weights mixture of experts model with 132B total parameters and 32k-token input context window, licensed for commercial use
Meta	LLaMA 3 (8B, 70B)	April 2024	Updated LLM with 8k-token input context window, with updated license allowing certain commercial uses
Mistral	Mixtral 8x22B	April 2024	Open-weights mixture of experts model with up to 141B parameters and 64k-input context window, licensed for commercial use
OpenAI	GPT-4o	May 2024	Multimodal LLM supporting voice-to-voice generation and 128k-token input context window
OpenAI	o1	September 2024	Reasoning model with built-in chain of thought (CoT) for complex scientific and mathematical problems
DeepSeek	R1	January 2025	Open source reasoning model (MIT license) optimized for math, coding, and logic

Today, you can build entirely new capabilities that would not have been possible previously to improve the lives of your employees and better serve your customers. For example, an agentic application<sup>2</sup> can now empower a financial crime investigator by automating the initial research for a case—a process that previously took days

---

<sup>2</sup> An agentic application is a software that uses an AI agent to perform tasks, make decisions, or interact with users with a defined level of autonomy. AI agents are discussed in more depth in [Chapter 2](#).

of manual effort. The system can synthesize internal records with external news and deliver a preliminary report in minutes, allowing the human expert to focus on critical judgment rather than tedious data collection.

This level of cognitive automation represents a new capability unlocked by modern LLMs (we will explore this specific Anti-Money Laundering [AML] Investigation Assistant in detail in [Chapter 7](#)). But you also have to keep up with rapid changes in the core technologies and use techniques that have not been fully proven. We are all now at the cutting edge.

This diversity of options among the technologies and techniques is truly a great thing. In fact, we are just scratching the surface for the potential uses of LLMs in the enterprise. It's easy to imagine a future where these technologies are generating massive amounts of value for the enterprise, automating mundane tasks, and making new products and services possible.

In this chapter, we will briefly introduce what an LLM Mesh is and then take an in-depth look at the many different types of LLMs that can be appropriate for use in the enterprise. We'll discuss different characteristics of models and how models are built, trained, and published, and how they run and perform.

After reading this chapter, you should be able to think about how you would want to use different models for different applications in your business. Given this multitude of models, you will see why an LLM Mesh architecture is going to be a key part of your AI strategy going forward.

## What Is an LLM Mesh?

As the state of the art of agentic application evolves, protocols have emerged to support the development of more sophisticated agents. Agent-to-Agent (A2A) defines how different autonomous agents communicate with one another, while a Model Context Protocol (MCP) standardizes how conversational history and state are packaged and sent to an LLM to ensure consistent behavior.

While these protocols help structure agentic logic, they don't solve the broader challenges of deploying a fleet of such applications in an enterprise. It is certainly possible to build the first few agentic applications without a formal architecture; in this common approach, the

logic for connecting to LLM services, managing access, and logging is built directly into each application. This is perfectly appropriate for a first proof of concept or a single application.

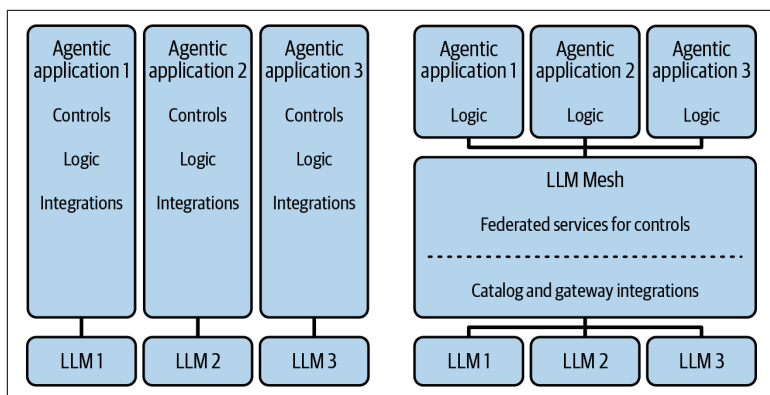
However, this direct, monolithic approach creates significant challenges at scale. As the number of applications, development teams, and overall complexity grows, an ad hoc approach leads to duplicated effort, inconsistent governance, and mounting maintenance costs. Addressing these enterprise-grade scaling challenges requires a new architectural paradigm.

The LLM Mesh is the architectural paradigm designed to solve these problems. There are three principles regarding what an LLM Mesh should accomplish. It should enable you to:

1. Access various LLM-related services through an abstraction layer (the importance of which is explained in [Chapter 2](#)).
2. Provide federated services for control and analysis (Chapters [3](#) through [6](#) describe these federated services).
3. Provide central discovery and documentation for LLM-related objects (the concept of “objects” in an LLM Mesh is developed further in [Chapter 2](#)).

These principles allow for agentic applications to be built in a modular manner, simplifying their development and maintenance. This approach directly borrows from foundational software design principles: each object in the Mesh is highly cohesive, focusing on a specific function, while the architecture ensures they are loosely coupled, allowing individual components to be replaced or upgraded with minimal impact on the rest of the system.

[Figure 1-1](#) illustrates monolithic agentic applications built outside of an LLM Mesh architecture and the same applications built within an LLM Mesh architecture. The LLM Mesh architecture creates a modular layer between the application logic and the controls for safety, security, and cost—and integrations with LLMs, enterprise data sources, and other systems.



*Figure 1-1. Comparing monolithic agentic applications built outside of an LLM Mesh architecture with the same applications built in an LLM Mesh architecture*

### Why LLMs and Not Generative AI?

An LLM Mesh architecture focuses on LLMs and not generative AI more broadly.

LLMs are large neural networks trained on text data. They possess a variety of natural language processing capabilities. Many, but not all, LLMs can generate text. Generative AI (GenAI) is a broader category of AI that includes models that can generate text, audio, images, and videos.

Beyond simply generating text, LLMs are also used to reason through a problem, to give instructions to various tools, and to write the code to connect to various tools. While image-generating models, for example, can be useful in the enterprise, they are not relevant in the context of building sophisticated AI applications that are the focus of the LLM Mesh.

An LLM Mesh provides a gateway not only to LLMs, but also to the full range of objects that are needed to build fully featured, agentic applications. These objects include the LLMs themselves and the services to host them, but also agents, tools, retriever services, and applications such as chatbots.

These objects are, for most organizations, new types of assets that will need to be developed and used. The skills to develop and use these kinds of objects are not yet commonplace in organizations,

and best practices for their development and use are still being defined. Amid this rapid innovation, the LLM Mesh architecture paradigm aims to simplify the management and use of these objects to accelerate and standardize the development of agentic applications ([Chapter 2](#) will explore in depth these different types of objects and how an LLM Mesh can simplify their use).

## The Right Model for the Right Application

The challenge for the use of LLMs in the enterprise is not a lack of availability of models, but selecting the right one for the task. As of December 2025, the popular model repository Hugging Face lists [2,254,112 models](#), of which 299,171 are [text-generation models](#). More models are being developed and released every day.

In fact, the abundance can actually be a hindrance, as you have to sort through the many different options to choose the ones that are best for your applications.

A large general model that can do most things pretty well is a good place to start. But as an enterprise's use of LLMs matures and it seeks solutions to niche problems, higher levels of performance, and optimized budgets, it will need to use a growing number of models across different applications.

The following sections explore the different practical considerations of models and how these characteristics may make a model more or less appropriate for the many different, specific uses in the enterprise.

## Model Size: The Upside and Downside of More Parameters

The word *large* in *large language model* refers to two interconnected concepts: the number of parameters in the model and the immense size of its training data, measured in tokens. While not a direct causal link, research has shown that the two are closely related; successfully training a larger foundational model requires a proportionally larger volume of training tokens to ensure high performance and prevent undertraining.

LLMs often have hundreds of billions to trillions of parameters. For example, GPT-3, released in May 2020<sup>3</sup> and the immediate precursor to the model behind the first version of ChatGPT, has 175 billion parameters. The first version of LLaMA from Meta AI in February 2023 had 65 billion parameters.<sup>4</sup> Increasingly, as the market becomes more competitive, the makers of proprietary models are no longer making the number of parameters in their models public.

These parameters are the numerical values (sometimes they will be called weights and biases) that make up the simple mathematical formulae of each neuron in the neural network. Usually, they are 32-bit floating-point numbers. A process called *quantization* can simplify these numbers to 4- or 8-bit integers. This process can often dramatically improve the efficiency of a model while having only a modest impact on model performance. This simplification process is illustrated in [Table 1-2](#).

*Table 1-2. Example of a single parameter value at different precisions*

Precision	Example value	Description
32-bit floating-point (FP32)	0.859375	The original, high-precision value. It requires more memory and computational power.
8-bit integer (INT8)	109	A quantized version. The value is mapped to a range of 256 possible integers (e.g., -128 to 127), saving significant space.
4-bit integer (INT4)	6	A more aggressively quantized version. The value is mapped to a tiny range of just 16 possible integers (e.g., -8 to 7), offering maximum efficiency.

[Figure 1-2](#) illustrates a simple neural network architecture, showing the input layer, two hidden layers, and the output layer. The circles represent the nodes in the network, the values under the nodes are the biases, while the values on the lines connecting the nodes represent the weights. In simple terms, weights determine the strength of the connection between neurons, while biases fine-tune the neuron's output. Larger neural networks, like LLMs, are built on the same

3 Tom B. Brown et al., “Language Models Are Few-Shot Learners,” preprint, arXiv, May 28, 2020, <https://oreil.ly/9g2Ti>.

4 “Introducing LLaMA: A Foundational, 56-Billion-Parameter Large Language Model,” Meta Research, February 24, 2023, <https://oreil.ly/3jRv3>.

basic architecture but are billions of times larger with more than one hundred hidden layers.

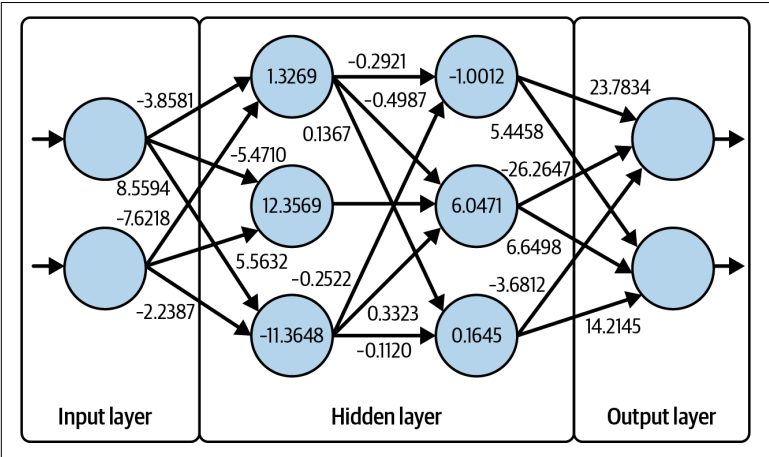


Figure 1-2. Simplified example of a neural network showing the input, hidden, and output layers and the weights connecting each node and the biases of each node

Generally speaking, larger models perform better on standard benchmarks for model performance. Thus, it could be easy to conclude that you should choose the largest model your budget allows and use it for everything. But that would be like using your large, comfortable, powerful grand touring car for every trip. While it would be the right choice for a cross-country road trip, it would be overkill for a quick trip to the grocery store or the bakery around the corner. A bicycle or your own two feet would be better for such errands.

The following subsections explain the trade-offs related to the size of a model.

### Inference costs

The most direct impact of a larger model size will be on inference cost. Inference is the process of generating tokens in response to a particular input. A model with more parameters will require more calculations during inference. Ultimately, these calculations run on physical hardware that consumes electricity, creating tangible infrastructure and energy costs.

In some cases, companies offering these models as a service may obfuscate these costs, for example, by subsidizing the cost in order to gain more customers. This may make an apples-to-apples comparison difficult. We'll dig into cost considerations in [Chapter 3](#).

Some models function, in essence, as a combination of smaller “expert” models that learn to specialize in processing different types of information. This architecture, known as mixture of experts (MoE), can dramatically reduce the cost of inference. One well-known model using an MoE architecture is 8x7B (also known as Mixtral) from Mistral. Despite being a 46.7-billion-parameter model, only 12.9-billion parameters are used per token. This is achieved by a small “router” network that directs each token to only the most relevant handful of “experts” for processing. This approach has led to improvements in inference cost, but makes the model more challenging to build and to fine-tune.

All other things being equal, larger models will be more expensive to use, though technological advancements like MoE mean that these trade-offs will become more complex in the future. The benefit that large models bring to a particular use case may justify their expense in certain cases, but a wise strategy will use them only where needed. The costs of different models hosted in different ways will be discussed in [Chapter 3](#).

## Inference speed

While the inference of a larger model will require more calculations, these calculations can be done more quickly when using larger- and higher-performance hardware. Furthermore, the overall inference speed is also a product of the entire system, influenced by factors like network bandwidth, request batching, and optimization techniques such as model quantization.

Standard benchmarks are being established to accurately quantify and compare the speed of different models, acknowledging that hardware and network performance will have a significant impact on the results. The two metrics that are used most commonly are latency and throughput:

### *Latency*

This is often measured as *time to first token* (TTFT), a measure of how long the model takes to generate its first response token to a user's input. In applications where the end user is interacting with

the model in real time, latency will influence whether the model “feels” responsive. In applications where the model’s response is part of a longer chain of interactions, latency will need to be considered when setting when the application will time out.

### *Throughput*

This is often measured as tokens per second (TPS), the overall rate at which the model will generate tokens in response to a given request. Like latency, it will influence if a model feels fast to an end user. Throughput needs to be taken into consideration when building applications that depend on the output of the model.

When comparing the speed of models, pay close attention to the units being used; different testers are using different methodologies.

While the relationship between model size and inference speed is indirect (because large models can be run more quickly on higher-performance hardware, and factors like network performance can influence the time it takes to receive a response), the measured speed of a deployed model must be taken into account when building an agentic application.

### **Task coverage and performance**

One of the main functional differences between LLMs and previous generations of models used for natural language processing (NLP) is that those earlier models were always task specific. For example, separate models would be used for sentiment analysis, text summarization, or language translation.

LLMs can do all of those tasks, and generative LLMs can do something that previous models could not: generate content based on a prompt. Generally speaking, models with more parameters can perform more tasks, which can be useful in the enterprise when several tasks need to be performed on the input text.

As LLMs scale in size and training data, they often exhibit surprising new capabilities. While initially described as “emergent abilities” that appeared unpredictably, the nature of this phenomenon is a subject of ongoing research. Some studies suggest these sharp performance gains may be an artifact of how performance is measured, rather than a fundamental shift in the model’s reasoning. Regardless of the cause, it is clear that larger models are often capable

of handling a wider and more complex range of tasks than their smaller predecessors.

In addition to gaining new abilities as they grow, larger LLMs generally show better performance on any existing task that they are capable of performing as well. Recent research shows that this improvement in performance is not linear nor predictable, as with the emergent abilities mentioned previously.<sup>5</sup>

## Context windows

The amount of input text that a model can receive within a single prompt is known as its *context window*. Measured in tokens, it defines how much information a model can work with at a single time.

For example, a model with a small context window can only be used to summarize a document that can fit in its context window. You could break up the document into smaller pieces, but the model would summarize each separately, without knowledge of the entire document, potentially resulting in repetitive or incoherent results. Large context windows, on the other hand, allow for plenty of space to provide examples of what you want the LLM to produce (called *few-shot learning*) and to engage in more complex prompt engineering techniques.

Generally speaking, larger models have larger context windows, and some models have been optimized for exceptionally large context windows. While the original GPT model had a context window of only 512 tokens (approximately one page of text), Gemini 1.5 from Google now has a context window of more than 1 million tokens and has been shown in internal testing to handle up to 10 million tokens, roughly equivalent to 15,000 pages of text.

## Sizing models to the task

Reading the previous sections, it is easy to conclude that if cost and complexity are no barrier, then the largest models are always the best choice for any application in the enterprise. But in which enterprise are cost and complexity not a barrier? In fact, these

---

<sup>5</sup> Ethan Caballero et al., “Broken Neural Scaling Laws,” arXiv, July 24, 2023, <https://oreil.ly/DswBL>.

are the two greatest barriers to the practical use of LLMs in the enterprise!

Given this reality, enterprise users of LLMs will need to choose a model that strikes the right balance of ability, performance, cost, and complexity for a specific application. The right choice for one application may not be the right choice for another application.

## General Models Versus Specialized Models

Building on this understanding of the implications of model size, we will now explore the differences between general models and specialized models.

General models are those that have been trained to perform at human level across a wide range of tasks. OpenAI's **GPT-4o** (released in May 2024) is an example of such a model. It demonstrates very high performance across a great number of tasks, covering natural languages, programming languages, and a wide variety of specialized jargon. It can generate, summarize, and translate text; it can write technical reports, and it can write poetry. Furthermore, as a natively multimodal model, GPT-4o seamlessly processes text, image, and audio data as input.

In contrast to these general models, specialized models have been trained to perform well on specific tasks in specific domains or have been compressed and optimized for performance at a smaller model size.

Note that while high-performing, general models tend to be larger models, specialized models may be larger or smaller. Model performance is explored in **Chapter 4**.

### Types of specialized models

*Task-specific models* are those that are focused on doing specific tasks very well. Some examples of task-specific models include **M2M100**, a model that is designed to translate between any pair of natural languages, or a model that is specialized in summarization. These models are often smaller and more efficient than general-purpose models for their single, dedicated function.

*Domain-specific models* are those that are trained on the language of a specific domain. For example, OpenAI's **Codex** is a fine-tuned version of GPT-3 trained on a massive dataset of public code,

making it highly capable at generating, explaining, and translating programming languages. Other examples include BioMedLM,<sup>6</sup> a 2.7-billion-parameter model trained on biomedical literature and thus well-adapted to answering questions about medical topics, and BloombergGPT,<sup>7</sup> a 50-billion-parameter model trained on a very large dataset of financial documents designed to serve the financial services industry.

*Resource-constrained models* are models that have been compressed through various techniques to maintain good performance in their desired tasks or across a wide range of tasks while being less resource intensive to run. An example is MobileBERT,<sup>8</sup> a distilled version of the popular BERT model designed to be run on mobile devices.

*Embedding models* transform text into numerical representations called embeddings or vectors. These embeddings capture the semantic meanings of the text and the relationships between the different parts of the text. A common application is retrieval-augmented generation (RAG) where a corpus of text (e.g., thousands of documents) is converted into embeddings and stored in a database called a vector store.

*Reranking models* are used to refine the initial list of results from a retrieval system, improving their relevance to the end user's query. This two-stage process allows for an optimal balance of speed and accuracy. While the initial retrieval is fast, the reranker performs a more computationally intensive but fine-grained analysis on a smaller set of candidate documents.

These models come in various forms. For example, models like the popular BAAI/bge-reranker-v2-m3 are highly efficient and specifically trained for this task. Additionally, there is active research into using general-purpose LLMs as rerankers, such as the Rank-R1 model, which is trained to reason about the query and documents

---

6 Elliot Bolton et al., “BioMedLM: A 2.7B Parameter Language Model Trained on Bio-medical Text,” preprint, arXiv, March 27, 2024, <https://oreil.ly/CWrjH>.

7 Shijie Wu et al., “BloombergGPT: A Large Language Model for Finance,” preprint, arXiv, December 21, 2023, <https://oreil.ly/mq9m2>.

8 Zhiqing Sun et al., “MobileBERT: A Compact Task-Agnostic BERT for Resource-Limited Devices,” preprint, arXiv, April 14, 2020, <https://oreil.ly/XuWNo>.

to improve its ranking capabilities. Each approach presents different trade-offs in terms of performance, cost, and complexity.

### Choosing a general or specialized model

The existence of a diverse and growing ecosystem of both general and specialized models gives enterprises the opportunity to use different models for different purposes.

In the enterprise, general models are well-suited to tasks where the input is going to be highly unpredictable. This could be the classification of documents into different categories. For example, if a directory contained a mix of contracts, invoices, and emails, a first step in the analysis could be to use a general model to sort the documents into different categories so that the contracts could be analyzed separately from the invoices.

Specialized models are well adapted for tasks where the input data is more homogeneous and predictable. Let's explore what this might look like across a hypothetical pharmaceutical company. That company may wish to build a chatbot to serve its customers (doctors, nurses, pharmacists, and other healthcare providers) in their interactions with patients. It would likely choose a domain-specific model like BioMedLM to ensure higher quality and more relevant results. The same company may then use a model like **ESM Meta-genomic Atlas** from Meta AI researchers, which has been trained on the language of proteins as part of its molecular research applications. Finally, that same organization may use a non-LLM computer vision model to watch its products as they come off of the manufacturing line to quickly identify any anomalies as part of its quality assurance processes.

General models can be a very good starting point for enterprises as they experiment and build their first use cases using LLMs. At those early stages, the simplicity of using a single model for a variety of tasks and use cases outweighs the benefits of further optimization using specialized models. But, as an enterprise scales its use of LLMs across use cases, enterprises will want to optimize their use to improve performance and reduce costs. In this context, specialized models become more relevant, and the number of models that an organization will need to manage and apply will tend to increase.

## What Is Fine-Tuning?

A common way of creating a domain-specific model is to fine-tune an existing base model. For example, Google created Med-PaLM 2 by taking its general PaLM 2 model and fine-tuning it on a curated dataset of medical knowledge, significantly improving its performance on medical questions and reasoning.

Fine-tuning is a type of transfer learning that feeds previously unseen—usually specialized—data into a model to retrain some parts of the model. Compared to building a model from scratch, it is far less complex and compute-intensive.

While fine-tuning is simpler and less expensive than building a base model, it remains an advanced technique and should be used only when other, simpler, and less expensive avenues have been exhausted.

Fine-tuning has often been cited as a way to elicit better performance from base models, allowing enterprises to differentiate their use of LLMs from their competition. While this is true, it ignores or downplays the difficulties of fine-tuning—such as the need for expert data curation, significant compute costs, and the risk of degrading the model's general abilities—and leaves unexplored the opportunity to generate differentiated results using simpler techniques like prompt engineering and RAG.

## Making Sense of Model Licenses

There has often been a conflation between a model's license (e.g., open source versus proprietary) and where the model is hosted (e.g., provided as a service via API versus self-hosted or on premises). It is important to distinguish between the two dimensions. For example, hosted services like Amazon Bedrock serve both proprietary and open models, while providers like Cohere license their proprietary models for self-hosting in addition to hosting the model themselves. Hosting options will be covered in the next main section, while this section will distinguish between the different license types.

### Proprietary models

Proprietary models are just that: proprietary to their creators. The creator of a proprietary model retains full control over the intellectual property of the model itself. Most often, these models are a

black box. In other words, their training data, the algorithm used to train the model, any subsequent steps such as reinforcement learning or fine-tuning, and the weights of the model itself remain hidden from the end user, unless the developer chooses to disclose any of this information.

Early in the development of LLMs, there was a trend toward openness, even among developers of proprietary models. The release of subsequent models has not been accompanied by such detail.

The use of proprietary models is governed by the terms of use that a customer agrees to when using the model. An enterprise should ensure a full and detailed legal review of these terms to ensure that they are appropriate for the intended use. Specific attention should be given to any rights that the model provider may claim to have on any data sent to the model for inference. Generally, models that are licensed for professional use do not retain any customer data for retraining purposes, though they may retain customer data for quality assurance purposes.

### **Open-weights models**

An open-weights model provides public access to the pretrained parameters of the model. This can allow the end user to modify the weights through fine-tuning or other techniques to adapt the model to their needs.

Open-weights models typically do not publish their training data, training algorithms, or other associated information. Meta's Llama 3 models are a prominent example; their weights are publicly accessible for anyone to use and modify, but the massive, curated dataset they were trained on is not released. As such, it can limit the ability to perform a detailed technical inspection of the model or to reproduce the model's performance. These limitations, however, are most relevant to other researchers and are less relevant to enterprises that are seeking to simply use a model in the most efficient and effective way.

### **Open access models**

Open access is a growing category of models that are nearly open, but have custom terms that cannot be considered fully open source in the traditional definition of that term. It covers a wide gamut of licenses with different restrictions, and thus should be the subject

of a detailed legal review to ensure that the license allows for the intended use.

Some examples include:

- BLOOM, which was released under the OpenRAIL-M license.<sup>9</sup> Though quite nearly open source, it has requirements for responsible use of the model, which means that it is not fully open source.
- Llama 2 and 3 from Meta AI have been released with their own custom licenses (called the **Llama 2** and **Llama 3** Community Licenses, respectively) that set limits to the use of the model. Specifically, the licenses forbid the use of the model in applications with more than 700 million monthly active users and for the purpose of building competitive models.

## Open source models

Open source models are the most open of all, publishing details of their training data, training algorithms, and model parameters, allowing for the most permissive use of the model. Common open source licenses include Apache 2.0 and MIT. Meta AI's first version of its Llama model was released under a restrictive, noncommercial license focused on research use cases, which made it unsuitable for most enterprise applications.<sup>10</sup>

## Choosing a license for enterprise use

Even though proprietary models are the most restrictive, they are often entirely appropriate for use in the enterprise, as with any other proprietary enterprise software. By charging for access to their models, providers of proprietary models may be able to more easily provide for services and support for the use of their model. This may make them more appropriate for use in the enterprise.

Open-weights, open access, and open source models may be more useful in applications where an enterprise wants more control over

---

<sup>9</sup> *BigScience* (blog), "The BigScience RAIL License," by Francesca Rossi et al. Posted May 20, 2022, <https://oreil.ly/Sdm1k>.

<sup>10</sup> Hugo Touvron et al., "LLaMA: Open and Efficient Foundation Language Models," preprint, arXiv, February 27, 2023, <https://oreil.ly/Exfat>.

the model itself and possesses the technical expertise to make any such modifications or to host the model.

## Model Hosting

Enterprises have three main hosting options when looking to access LLMs:

- API services from the model developers, such as OpenAI, Anthropic, Cohere, and Mistral.
- Cloud service providers offering hosted LLM services, such as Azure OpenAI Service, Amazon Bedrock, or Google Vertex AI Model Garden. These services also allow customers to load their own models, while the underlying hardware is managed by the cloud provider.
- Self-managed hosting of models. Many models with different licensing terms are available for self-hosting, including the open source and open access models described previously. Cohere also licenses its proprietary models for self-managed hosting.

In many cases, models hosted by their developer or a cloud service provider are the best choice in the enterprise. In the same way that cloud computing outsourced the burden of running data centers, hosted models are a simple continuation of that trend, offering infrastructure as a service (IaaS). Given the intensive compute requirements of LLMs, especially under heavy workloads, outsourcing this can be a wise choice.

The most common objection to using a hosted service is that it requires sending corporate data to a third-party service. But in many cases, this corporate data is already hosted by a third party that may also be hosting internal communications and other sensitive data (e.g., a company that uses Microsoft 365 productivity and communication tools has its data in Azure). Is using the LLM service from that same provider any different? It is ultimately a question that warrants review by your legal and risk teams.

Self-hosting a model requires acquiring the necessary hardware, configuring it to run the LLM, and then maintaining that stack for reliable internal use. For larger models, this will require a cluster of GPUs that have been properly configured with the right drivers and

packages to run the LLM in question. The LLM must then be loaded into this environment so that it can begin to serve internal requests.

Self-hosting can be an appropriate choice for an enterprise in cases where an organization needs full control over the model and the hardware it runs on and cannot use a third-party service for its data. This may be the case in the most restrictive data environments, or if the enterprise does not want to rely on a third party to ensure the performance of the environment, notably in contexts where third-party providers may need to throttle access to certain customers to ensure the overall stability and availability of their service.

In the case of both self-hosting and hosted services, applications that use the LLM will usually access the model through an API endpoint. The difference is simply who is hosting and maintaining that endpoint and whether the data going to and returning from the LLM leaves the corporate firewall of the enterprise.

## **Building a Base Model Is Not for Most Organizations**

Early on in the popular interest in LLMs, a lot of attention was given to the expense and complexity of building these models. Billions of dollars were being spent building these models, and sometimes training them took many months. A huge amount of this initial work was amassing the enormous training sets required to build models of this scale.

Recent advances have brought down the time needed to build new models, and open source training data repositories now exist. But the fundamental question for an enterprise that is considering building a model remains: why would you? Given the great diversity of today's models, which offer seemingly endless combinations of performance, specificity, licensing, and hosting options, what would justify the time and expense needed to build your own model, especially given that you are uncertain of being successful?

Any company whose core business is not building or serving AI models should not consider building its own model. There are more than enough options on the market today. The challenge is not getting access to a model but using it safely, securely, efficiently, and effectively to further your business goals. This is where an LLM Mesh comes into play.

## Bottom Line: Why the LLM Mesh?

As you have read in the previous sections, a great variety of models exists in an ecosystem that is rapidly evolving. This is ultimately a very good thing for enterprises: it means that they will be able to pick and choose the right model for the right applications within their business. Building applications that are powered by these LLMs requires combining them with other objects, like retrieval systems, prompts, and tools. This requires careful attention to many different factors:

- How the models, services, and associated objects are registered and used within the organization
- How the data is routed to the model
- How access to the models and services is controlled
- How the use of the model is logged and audited
- How the content generated by the model is moderated
- How the models can be enriched with proprietary data
- How the applications can be developed, deployed, and maintained efficiently
- How more people can become involved in this process

As more agentic applications are built and used in the enterprise, the cost and complexity of managing all of these dimensions risks spiraling out of control. This could force the enterprise to make compromises, potentially limiting the value that it derives from AI.

For example, consider a company that has a code assistant that is well-versed in the company's proprietary code libraries. This is a use case that would benefit from using a small, specialized model that is self-hosted and to which access is restricted. But if the organization lacks the ability to quickly and efficiently add this model to its mix, it may not pursue this use case, leaving the potential gains in efficiency on the table and falling behind its competition.

This code assistant could have provided a real efficiency gain to the developers working at the company. Missing out on such an improvement would be unfortunate, given that many of the additional capabilities that are required to use an LLM efficiently and effectively in an enterprise are common to all models.

This is the power of an LLM Mesh: its ability to reduce the cost of building an additional agentic application in the enterprise. With an LLM Mesh, an enterprise is free to develop an optimal AI strategy without compromising on performance, cost, safety, or security.

The remaining chapters of this technical guide will go into much more detail about how implementing an LLM Mesh can be done.



# Objects for Building Agentic Applications

LLM Mesh is a new architecture paradigm for building agentic applications in the enterprise. It enables an organization to build and maintain more agentic applications, ultimately getting more value from LLMs.

An LLM Mesh will allow you to:

- Access various LLM-related services through an abstraction layer
- Provide federated services for control and analysis
- Provide central discovery and documentation for LLM-related objects via a catalog

This chapter will describe the many different types of LLM-related services that are used in building agentic applications and how they can be connected with one another. These various LLM-related services are called *objects* in an LLM Mesh. The final section of this chapter will describe the importance of the catalog for the discovery and use of the objects in an LLM Mesh.

We start with an explanation of why using an LLM Mesh to build agentic applications is increasingly important in today's competitive landscape. The bottom line is that you are going to need to build a lot of custom agentic applications.

# The Potential of New Agentic Applications

In [Chapter 1](#), we learned about the many different types of models available, how they work, and the options available for hosting them. What can these models be used for in the enterprise? There have been two main, initial uses of these models in the enterprise. The first is simply providing a version of the consumer chatbot experience within a wrapper that meets enterprise security and auditability requirements. The second is using these models to provide software assistants, often called *copilots*, that can accelerate the use of existing SaaS products.

This first generation of enterprise use has been met with a mixed reaction. In some cases, notably when used as coding assistants for software developers, the copilots have proven to be valuable additions to the enterprise IT mix. Other feedback has been more varied, such as a study by MIT showing that 95% of GenAI pilots do not make it to production, leading in some cases to disillusionment. It is too early, however, to discount the potential for LLMs in the enterprise. This is because the second generation of agentic applications in the enterprise will be more capable and more valuable.

These applications will not only use the ability of these models to generate text but also their ability to solve arbitrary problems when instructed to do so. For example, an LLM could be provided with the documentation for an API that looks up the current price of a stock. With that documentation and its coding ability, the LLM can write a script to call that API for a given stock price. If allowed to execute that script, the LLM—without having ever been explicitly programmed to do so—could then become a tool for the end user to look up arbitrary stock prices. This ability to accomplish tasks for which the LLM has not specifically been programmed is called *generalization*, and it is what allows LLMs to be the engines in a new class of enterprise applications by handling the unexpected and ambiguous contexts that limit traditional applications.

These new applications will provide automation and decision support throughout the enterprise. In order to do so in a reliable and cost-effective manner, however, they will need to be carefully designed, tested, deployed, and monitored. While many of the constraints of traditional enterprise applications, such as security and reliability, will also apply to this new class of agentic applications,

the way in which they are built and the components that are used to build them will be different.

Given the LLM's ability to generalize, would it be possible to develop a single, all-powerful application that can solve any problem and answer any question in the enterprise? In short, no. While the LLM itself is capable of generalization, the constraints of the enterprise will require that the scope of any one application be relatively narrow to ensure consistently good performance and to control access to data and tools.

For example, this imaginary, all-powerful application sounds convenient but would require full access to all of the company's data and tools, from the most mundane to the most sensitive. Just as an employee should only have access to the data and the tools that they need to do their job, so too must the access of any one agentic application be limited to that which it needs to perform its function. Furthermore, while LLMs are capable of generalization, they require quite specific instructions to deliver consistent results, often with examples of the expected input and output. This also drives toward a larger number of more narrowly scoped applications.

Concretely, how many such agentic applications might a large corporation need? Let's do some order-of-magnitude estimations. Let's say that a large corporation has 10 departments and each department has five core functions. Each of these functions could potentially benefit from five such applications. For example, within a sales department, the sales operations function could have one application that researches the department's target accounts, a second that checks if the sales process is being respected, a third that continuously analyzes the health of the sales pipeline, a fourth that summarizes meetings with prospects, and a fifth that assists salespeople with their follow-ups.

Doing the multiplication of this order-of-magnitude estimate gives us  $10 \times 5 \times 5 = 250$  such applications in that enterprise, though some applications could be shared by departments. Again, this is not a precise number; it's a rough estimate of the order of magnitude of applications to expect. Expecting several hundred such applications in use in a large organization seems like a reasonable estimate.

# Build Versus Buy

If an organization would benefit from several hundred novel applications, where would they come from? As always, organizations will face a “build versus buy” decision. On the buy side of that balance, existing software vendors and new startups are already bringing these applications to market, and organizations will have a lively marketplace of competitive offers to choose from. On the build side of the balance, more advanced organizations are building their first production-ready agentic applications, such as a pharmaceutical company that is partially automating its **patent case law review process**.

Which approach is best? Each has its advantages, disadvantages, and appropriate uses, meaning that most organizations will buy some applications, and build others. **Table 2-1** summarizes these trade-offs and considerations.

*Table 2-1. Comparing the trade-offs of building versus buying agentic applications*

	Advantages	Disadvantages	Ideal uses
<b>Buying off-the-shelf agentic applications</b>	Turnkey performance once implemented Developed and maintained by professional software engineers	Same performance as your competitors that use the same solution Complex to integrate with enterprise systems Governance challenges for tracking which models are used by which applications	Noncritical functions where the goal is to gain in efficiency, not necessarily to differentiate from competitors
<b>Building custom agentic applications</b>	Adapted to specific business context with the potential to build differentiated capabilities Full control and transparency over the application Independence from software, AI, and cloud providers	Skills required to build applications may not be available Complexity of monitoring and maintaining grows with the number of applications	Core and strategic functions where full control and strong competitive differentiation are needed

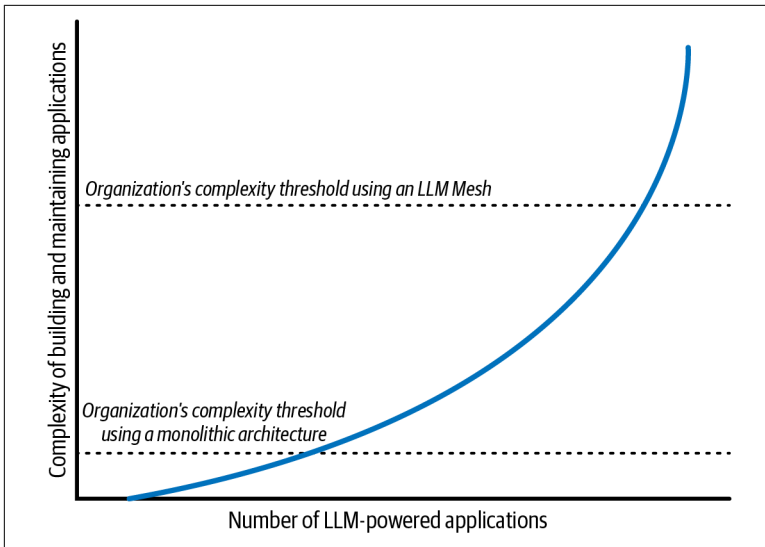
Agentic applications, whether they are custom-built or bought off-the-shelf, have the potential to improve the efficiency of an organization's operations. But simply improving your efficiency in lockstep with that of your competitors does not improve your competitive position in the market. If you are making the same efficiency gains as your competitors and no more, you are not becoming more competitive; you are simply keeping up.

Building custom agentic applications allows an organization to create a capability that its competitors do not possess and thus to outperform them in that particular domain. Given the cost and complexity of building, monitoring, and maintaining these applications, organizations will choose to focus their internal development efforts on the parts of their business that stand to benefit most from strong competitive differentiation. In most cases, this will be their core business. For example, it may be R&D and supply chain management for a pharmaceutical company or risk and price modeling for an insurance company. The needs of noncore functions will be satisfied with applications bought off the shelf.

## The Complexity Threshold

How many custom agentic applications can any given organization develop and maintain? This limit isn't just about the number of applications but the compounding complexity created by how they are built. Traditional development practices often create applications in silos. While this solves immediate problems, each new application multiplies the web of interdependencies and adds to the organization's technical debt. We call this the organization's *complexity threshold*, and it is illustrated in [Figure 2-1](#).

As the organization develops and deploys more agentic applications, the complexity of monitoring and maintaining them increases until, at some point, the maximum complexity is reached and no more applications can be developed. Reaching this threshold means that the organization cannot develop more applications, even if doing so would benefit its business. If the organization wants to develop more applications, it must find a way to increase its complexity threshold. This requires standardizing and structuring the way that the organization builds these applications.



*Figure 2-1. An organization's complexity threshold*

## A New Paradigm for Building Agentic Applications

Bringing standardization and structure to the way that applications are built in the enterprise is a show we've seen before. Over the years, organizations have used different architecture paradigms for developing applications. Starting with monolithic applications in the early days of application development, where all components were tightly integrated into a single codebase, organizations then shifted to an architecture paradigm with a higher degree of abstraction with the service-oriented architectures of the late nineties, and now, the modern standard of microservices has taken that abstraction even further.

Today, the dominant paradigm for building agentic applications relies on development frameworks like LangChain, AutoGen, CrewAI, and Haystack. While these toolkits provide modular, composable components, they are often used to build applications where the agent's core logic becomes highly centralized and tightly coupled.

This architectural pattern—where tools, prompts, and reasoning steps are interwoven within a single application—reflects the relative immaturity of agentic design in the enterprise. Even when built from modular parts, the resulting system often suffers from

challenges like difficult debugging, production fragility, and rigid, hard-to-maintain workflows.

A new architecture paradigm is needed for building and maintaining many agentic applications that can raise an organization's complexity threshold. LLM Mesh is that new architecture paradigm.

Now, let's look at the objects used in building an agentic application.

## LLM Mesh-Related Objects: An Overview

Building an LLM Mesh requires understanding the different types of objects that must interact with one another within an agentic application. [Chapter 1](#) covered the LLMs and the various services that host and serve them. While those models and services are at the heart of an agentic application, more is needed, especially if the developer hopes to build a custom application that will stand apart from the competition and deliver better and more valuable performance, as discussed in [“Build Versus Buy” on page 26](#). This requires integrating the LLMs with various objects unique to the organization.

An LLM Mesh thus treats objects of a similar type in the same way, with the LLM Mesh itself providing the translation between the generic object (e.g., a tool) and the specific service (e.g., a specific SQL database). This is critical from a maintenance and debugging perspective. For maintenance, a specific service like a database can be upgraded or swapped by updating its central connection in an LLM Mesh, without needing to change the code of every application that uses it. For debugging, this allows for the rapid isolation of problems—an engineer can quickly determine if an error lies in an application's logic or in the specific tool's implementation. In this way, we say that the LLM Mesh provides *abstraction* between the high-level object and the underlying, specific service.

[Figure 2-2](#) illustrates the objects of an LLM Mesh organized into different layers that comprise the typical stack of an agentic application, overlaying the typical stack of a traditional application.

Note that in [Figure 2-2](#), the objects in the white rectangles are not themselves part of an LLM Mesh, but rather are abstracted as the higher-level objects in the shaded rectangles. This will be discussed further in [“Retrieval Services” on page 37](#) and [“Tools” on page 45](#). In contrast, traditional applications use data querying services and API

services directly, without abstraction as tools. Unstructured data is not used directly in traditional applications but is first transformed into structured data using traditional NLP techniques.

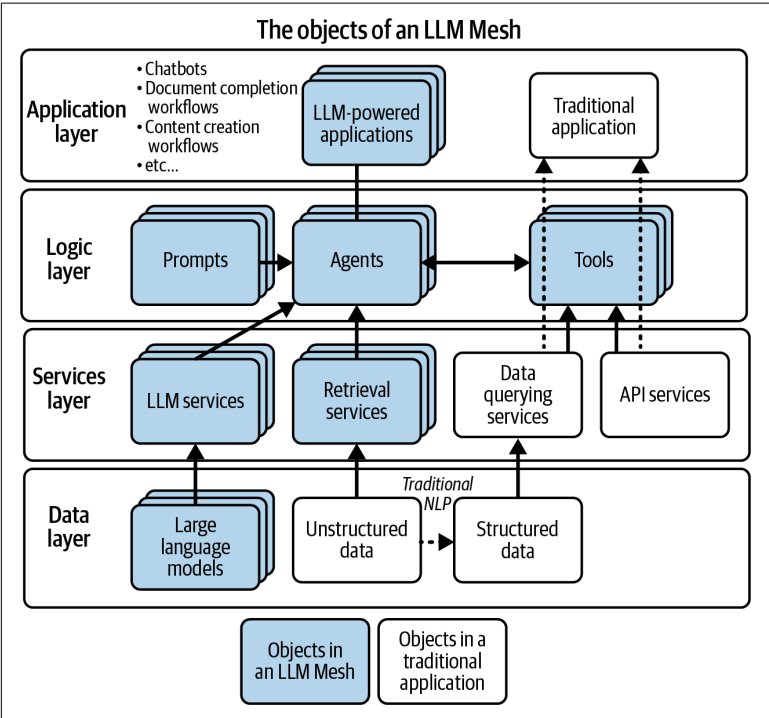


Figure 2-2. The objects of an LLM Mesh in comparison with those of a traditional application

Here is an overview of the different objects in [Figure 2-2](#) and how they relate to one another:

*Large language models*

The base model—the trained neural network comprising the core mathematical weights—as described in [Chapter 1](#).

*Unstructured data*

Enterprise data that is not in tabular form. A common type of unstructured data is documents, which may be in *.pdf*, *.docx*, or other formats. Unstructured data is abstracted as retrieval services in an LLM Mesh.

### *Structured data*

Enterprise data that is in tabular form, typically stored in databases, data warehouses, and data lakes. Structured data is stored in data querying services, which are in turn abstracted as tools in an LLM Mesh.

### *LLM services*

The services that are comprised of the hardware and software systems used to deploy and interact with the model in real time. As described in [Chapter 1](#), these services may be managed by the model developer, a third party, or internally by the enterprise.

### *Retrieval services*

A service that allows for the efficient and effective querying of unstructured data. The retrieval services usually consist of a specialized LLM used for embedding, storage for the embeddings (which can be either a dedicated vector store or another type of database—SQL or search, for example—that has added these capabilities), and some system for ranking the results to best respond to the query.<sup>1</sup>

### *Data querying services*

Databases and their associated query languages, like SQL, that allow for the efficient retrieval of structured data. These systems are abstracted in an LLM Mesh as tools.

### *API services*

Any internal or external API services to be integrated with the agentic application. An external example could be a weather service to look up a forecast, while an internal example could be the data catalog to allow for data discovery. These services can be very diverse and are abstracted in an LLM Mesh as tools.

### *Prompts*

The input to the LLM services can be templated and standardized and can run the gamut of prompting techniques (few-shot, chain of thoughts, etc., as described in the following sections).

---

<sup>1</sup> As retrieval services are relative newcomers to the enterprise architecture landscape and are themselves powered by LLMs, they are treated as distinct objects from other tools in an LLM Mesh.

### *Agent*

An LLM-powered system that seeks to accomplish a certain goal over multiple iterations within a defined level of autonomy and using tools to meet its objective. Note the centrality of agents in this architecture. They are the object where the logic and behavior of the application are defined.

### *Tool*

Any function or resource that an agent can use to accomplish its task. It can be a programming or querying language, an API service, or even another agent.

### *Agentic applications*

An application that provides a user interface and other functionality on top of the agent. A chatbot is one example of an application type, but agentic applications could have many different types of interfaces running the gamut from dashboards to mobile apps to assistants embedded in other applications to headless applications running behind the scenes and alerting users only when needed.

## **The Objects of an LLM Mesh in Detail**

The previous section provided a broad overview of all the components in a typical agentic application stack, including the underlying services. The following sections will now focus in detail on the seven core, high-level objects that formally constitute the LLM Mesh itself. The seven objects are LLMs, LLM services, retrieval services, prompts, agents, tools, and applications. Each section will first define and then describe the object.

At the end of each section, you will find a tip box titled “Thinking Like an LLM Mesh.” This box describes the expected input and output of each object. Recall that one of the main benefits of an LLM Mesh is that it creates an abstraction layer that standardizes the inputs and outputs of diverse services into standardized objects. The tip boxes summarize what those standardized inputs and outputs should be.

Building an agentic application requires integrating several different objects. For example, a simple chatbot application using a retrieval augmentation technique could be built using the following objects:

- An *application* with a chatbot interface where the end users ask their questions and receive their responses as well as provide feedback to the developers
- An *agent*, composed of several templated *prompts*, that defines how the user's question will be handled by the LLM
- An *LLM service* that receives the question, tokenizes it, and submits it to the *LLM* that will generate the response, enriching it with an answer from a retrieval service
- A *retrieval service* that provides access to unstructured data from documents; the retrieval service is comprised of:
  - An *embedding model* that converts the text data to vectors
  - A *reranking model* that will select the most relevant answer to the user's question, providing it back to the LLM service for inclusion in the reply

Such a chatbot could, of course, be built without an LLM Mesh simply by building a monolithic application that calls the various services, passing the results from one object onto the next. In practice, the developer of such an application would be writing many API calls, each of which is specific to each service. If the developer would later want to change, for example, from one third-party LLM service to another, this would require manually updating the code so that the application calls the new LLM service in the way that is expected by that service.

In that scenario, an efficient application design would provide a certain degree of abstraction, defining the interface with the LLM service as a single function within the application and not specifying the details of the API call in every instance where the LLM service is called.

An LLM Mesh takes this abstraction further, completely separating the service from the application and providing a standard interface for all objects of the same type for use across all agentic applications in the enterprise.

## LLMs

We covered LLMs in detail in [Chapter 1](#). When we talk about an LLM, we are talking about a very large file, often measuring in gigabytes or terabytes. For example, the 405 billion parameters of Meta's [Llama 3.1 model](#) weighs in at approximately 810 GB.<sup>2</sup> The majority of the data volume is taken up by the weights of the model itself. Remember, as described in [Chapter 1](#), the weights of a model are simply a great quantity of floating-point numbers.

If an organization is using a managed LLM service, they will never interact with the model itself, only with the service endpoint. But, if an organization self-hosts an LLM, then they will need to load the LLM into their hosting infrastructure.

### Thinking Like an LLM Mesh

From the perspective of an LLM Mesh, an LLM is thus an object that can be interacted with in only two very simple ways: it can be inferred and updated in the environment where it is hosted.

## LLM Services

An LLM service is a combination of storage resources, compute resources, and supporting software that allows an LLM to be hosted and accessed for inference.

The developer of the model may provide LLM services. For example, OpenAI, Anthropic, and Mistral all provide services that run their proprietary models. In these services, the end user does not load the model into the service; they simply select the service running the model that they prefer.

Alternatively, an organization may choose to build and run its own LLM service, managing the GPUs and associated technologies.

Finally, cloud service providers (CSPs) offer managed LLM services. In these, the end user may select the model they wish to run, but the CSP manages the compute and storage infrastructure.

---

<sup>2</sup> Parameter count and data type are from Meta's official announcement.

An LLM service, be it hosted by your organization or by a third party, is accessed via an API. Generally, most LLM services will expect similar variables when they are called. These include:

- Which model version to use
- A system prompt set by the developer to guide the model's completion
- The user prompt for the model to complete
- A temperature setting to define the level of randomness in the response
- Alternatives to temperature, such as `top_p` or `top_k`, which use different sampling methods to determine which subsequent token to select

In response to such requests, the LLM service will generally reply with a response that includes:

- An indication of the type of response (e.g., text completion or streaming chat)
- A unique identifier of the response
- The generated content
- Reasons for why text completion may have stopped
- Usage statistics about the number of tokens in the request and response

Most services have broadly similar expected inputs and outputs. An LLM Mesh abstracts and standardizes these inputs and outputs through its abstraction layer, ensuring that the request sent to a given service is formatted appropriately and uses the correct syntax. When using an LLM Mesh to build an application that calls an LLM service, the end user calls the LLM service object in the LLM Mesh, indicating which service to use, and the LLM Mesh translates that generic call into the specific call expected by the indicated service.

To better understand the value of providing a standard interface for all LLM services, let's compare the expected syntax of two common providers, OpenAI and Google Gemini, starting with OpenAI. The [OpenAI documentation](#) gives the following example:

```
curl https://api.openai.com/v1/chat/completions \
-H "Content-Type: application/json" \
-H "Authorization: Bearer $OPENAI_API_KEY" \
-d '{
  "model": "gpt-4o",
  "messages": [
    {
      "role": "system",
      "content": "You are a helpful assistant."
    },
    {
      "role": "user",
      "content": "Hello!"
    }
  ]
}'
```

Let's compare that with the expected request to the Google Gemini API. [Google's documentation](#) gives the following specification:

```
curl -X POST \
-H "Authorization: Bearer $(gcloud auth print-access-token)" \
-H "Content-Type: application/json" \
https://${LOCATION}-aiplatform.googleapis.com/v1/projects/ \
${PROJECT_ID}/locations/${LOCATION}/publishers/google/models/ \
${MODEL_ID}:streamGenerateContent \
-d '{
  "contents": [{
    "role": "user",
    "parts": [{
      "text": "TEXT"
    }]
  }]
}'
```

These short samples from the documentation already show some differences between the two APIs:

- OpenAI specifies the model in the JSON payload with the model key-value pair, while Google specifies the model in the URL path.
- The array containing the content of all messages is an object called `messages` by OpenAI and `contents` by Google.
- Google's format nests a `parts` array inside the `contents` array to hold the individual components of a multimodal prompt, such as combining a block of text with an image.

An LLM Mesh standardizes these and other differences, allowing for faster application development and easy switching between LLM services. As LLM service providers update their services, the LLM Mesh developer will update the Mesh accordingly, freeing the application developers from the need to do so.

### Thinking Like an LLM Mesh

As an object in an LLM Mesh, an LLM service expects a prompt as input and is expected to provide text as its output. This text can be natural language, function calls, tool invocations, or multiturn response instructions.

## Retrieval Services

From the perspective of an LLM Mesh, a retrieval service takes a user's query as its input and provides a relevant result from unstructured data as its output. It is how an agentic application accesses unstructured data. In this context, the unstructured data is text data coming from documents, often stored as PDFs, plain-text documents, or other common document formats, like *.docx*. Retrieval services allow agentic applications to make this data available to the employees of an enterprise by allowing them to discover it more accurately and rapidly.

Importantly, the information in these documents is not only made available to the employees. It will also be available for agentic applications themselves to inform their inference on how to solve a problem. Retrieval services serve this dual purpose: making the unstructured text data available to both the employee and the agentic applications, which, in both cases, leads to better decisions.

In retrieval services, like traditional search systems that came before them, there is usually a trade-off between the speed of the results and the quality of the results. While it is possible to have fast results or good-quality results, it is difficult to have both. Retrieval services are usually made of three separate components to provide the highest quality results as quickly as possible. These are:

### *Embedding models*

To convert the text of the document base, as well as the text of the query, into dense vector representations

### *Vector store, or other data storage*

For storing and indexing the vectors, allowing for efficient similarity search

### *Reranking models*

To improve the quality of the search results

Increasingly, these services are being provided as bundled services from various providers, as their combined functionality is required to provide the desired result to the end user: a relevant result from unstructured data in response to a natural language query. From the perspective of an LLM Mesh, the mutual dependency of these three underlying technologies is why they are combined as a single object, “retrieval services.”

The following sections describe the components of a retrieval service and some of the trade-offs that the various choices will entail.

## **Embedding models**

As described in **Chapter 1**, embedding models transform text into numerical representations called *embeddings*, stored as high-dimensional vectors. For example, the word *banana* might have the following embedding: [0.534, 0.312, -0.123, 0.874, -0.567, ...]. Each number represents the value of a particular dimension. If the model generates 100-dimensional embeddings, the complete vector for *banana* would be a list containing 100 numbers.

These embeddings capture the semantic meanings of the text, such that *Denver* and *capital of Colorado* will have similar vector representations (even though they share no keywords), while *kid* meaning *young goat* will have a different vector representation than *kid* meaning *young human*.

Different embedding models use different embedding lengths, meaning more or fewer dimensions for each vector. Put more simply, a shorter embedding length means fewer dimensions in each vector, and thus fewer numbers in the list that represents each word or part of a word. More embeddings require more storage and compute resources. New embedding models, like OpenAI’s text-embedding-3 family of models, allow for the embedding size to be shortened to a degree specified by the user. Shorter embeddings reduce storage and computational costs but may result in degraded performance, particularly for tasks requiring nuanced

semantic understanding. The actual impact on performance is highly dependent on the specific application and data. Model developers are working to increase the performance of vectors with fewer embeddings.

Embedding models expect input text that has been preprocessed to a certain degree. Different models have different requirements; an LLM Mesh provides a standard interface that is mapped to each model. Preprocessing will generally include extracting the text from any documents (e.g., *.pdf* or *.docx* formats), adding special tokens to tell the model about relevant breaks in the text, and splitting longer documents into smaller “chunks” that are sized appropriately for the embedding model.

In a retrieval service, embedding models serve the dual purpose of converting the corpus of documents into vectors and then doing the same for the query. Converting the corpus of text into vectors is usually an offline task, done once, while converting the query is necessarily done at runtime, when the query is received from the user.

### **Vector store, or other data storage**

The embeddings are written to a data store, often a dedicated vector store. A vector store is a database specially designed to store and efficiently query high-dimension, dense vectors, like those created by embedding models. Vector stores have built-in retrieval functionality for finding a stored vector most similar to the query vector, using functions such as cosine similarity, dot product, or Euclidean distance.

Traditional data stores—including relational databases like PostgreSQL, document databases like MongoDB, search engines like Elasticsearch, and graph databases such as Neo4j—are all adding support for dense vector data types. As traditional data stores add support for vector data, they are becoming a viable alternative to dedicated vector stores for some applications. However, these integrated solutions may not yet match the performance of specialized vector databases, a trade-off developers must consider.

This evolving technology landscape is one more reason why abstracting these services as “retrieval services” is important in an LLM Mesh. While the underlying technologies may change, the

function remains the same: provide relevant results from unstructured data to user queries.

## Reranking models

The vector store's retrieval function will provide a fast result, but it may not always be the most accurate. More accurate results can be obtained by using the vector store's retrieval function to narrow down the results and then a reranking model to analyze the subset more carefully, selecting the best result to return.

In contrast to the retrieval function of the vector store, the reranking model will take the entire source document plus the input query for comparison. Given that the source data could contain thousands or even millions of source documents, it would be too slow and too costly to run that process across every document. By using the retrieval function to narrow down the results to the top few (a number that can be specified) and then running the reranking model across the subset, you strike the best balance between speed and quality. This is called *two-stage retrieval*.

The retrieval system will provide the top-ranked results back to the agent. Often, multiple results are returned (the *top-k*) so that the LLM can cite and synthesize information from several sources.

### Thinking Like an LLM Mesh

As an object in an LLM Mesh, a retrieval service expects a natural language query as its input and is expected to output the top-ranked result from unstructured data. These results are then generally passed on to an agent.

## Prompts

Since the popular use of LLM-powered chatbots increased dramatically following the release of ChatGPT and associated products, many of us are now familiar with the notion of a prompt. A prompt is the initial input (a question, a command, or instructions) provided to the model, prompting its response.

In contrast to the ad hoc prompting often used in consumer applications, prompting in the enterprise benefits from a structured, templated, composable approach. This allows a bank of prompts to be developed, tested, and then shared for reuse across the organization.

While prompts can be used by any component that calls an LLM service, they are especially critical in defining agentic behavior. In this context, a sophisticated prompt or chain of prompts acts as the instructional code for an agent, setting its goals, constraints, and the logic it should follow. There are many different types of such prompts. The following sections discuss several categories of prompts, with some simple examples of each.

### **Role-based prompts**

These prompts direct the LLM to respond as a specific type of expert, such as a customer support agent or HR consultant, or guide the AI on the tone, formality, and style of its responses. Here are two examples:

“You are an IT support technician. Assist the user in troubleshooting their software issue.”

“Respond in a professional and concise manner suitable for senior management.”

### **Compliance and ethical prompts**

These prompts direct the LLM to provide responses that adhere to specific regulations or legal frameworks, or responses that follow specific internal guidelines for ethical practices:

“Ensure that no response contains personally identifiable information (PII), such as names, phone numbers, or identifiers like Social Security numbers.”

“Generate responses that respect the following internal ethical AI guidelines [corporate ethical AI guidelines].”

While using such prompt components can decrease the risk of non-compliance, they cannot guarantee that any result will necessarily be compliant. To reinforce human oversight, organizations can implement programmatic guardrails. These automated safety checks are explored in detail in Chapters 5 and 6.

### **Customization, personalization, and context-specific prompts**

These prompts customize responses from an LLM based on known information about a user or customer, a prediction about them, or other relevant contextual information. The variables in the example prompts here would be completed based on information in the

enterprise customer relationship management (CRM) system, customer support records, or using the result of a predictive model:

“Personalize the marketing message for a [age]-year-old [gender] living in [postal code].”

“Recommend to the customer [result from next-best offer prediction].”

“Given the customer’s previous request about [subject of the previous request], provide a relevant response.”

## Multistep process prompts

These prompts guide the LLM to respond in a multistep process by breaking down complex decisions into smaller, more manageable steps. These multistep process prompts are the building blocks of agents:

“Step 1: Gather all financial data from Q1. Step 2: Generate a financial report. Step 3: Summarize the key findings in a presentation.”

“First, evaluate the market demand. Next, assess the cost implications. Finally, recommend a go/no go decision.”

### Thinking Like an LLM Mesh

For the purposes of an LLM Mesh, prompts need to be tested, approved, and published in the catalog, which we will explain later in this chapter. Any prompt must be associated to a specific model and version, as small changes in the model or prompt may result in dramatically different prompt performance. These prompts can then be combined with one another to compose more complex and sophisticated prompts, themselves part of agentic applications.

## Agents

While various definitions for *agent* exist, from the perspective of an LLM Mesh, an agent is an LLM-powered system capable of accomplishing its objective across multiple steps using tools, without requiring prompting by an end user for each step. In an agentic application, the agent acts as the central orchestrator: it consumes prompts for its instructions, uses tools to interact with other systems, and chains these actions together to achieve its goal.

Within an LLM Mesh, an agent is the object where the other objects interact with one another to form a system that can respond to users' needs. They call one or more LLM services, they use several templated prompts, and they use one or more tools. As such, agents are some of the most important objects within an LLM Mesh and are at the core of building agentic applications in the enterprise.

Like the other objects, they must be built, described, cataloged, and maintained. As the maturity of an organization increases, it will begin to develop more agents and will likely start chaining those agents together, with one agent using another as a tool. This increasing complexity can be tamed by the abstraction and modularity that the LLM Mesh offers. Because each agent is treated as a standardized, reusable object, their interactions can be orchestrated without requiring custom integrations. The following chapters will explain how the Mesh's federated services make these complex chains manageable.

There are a few important parts to the definition of an *agent*, so let's look at them one by one. All of these parts would be defined in the agent's prompt.

## Objective

An agent's developer will define its objective by giving it a role-based prompt, as described in the previous section. For example, an agent that is part of an application that is designed to generate real-time sales analytics could include the following role-based prompt template:

You are a Business Intelligence Analyst with access to the company's sales data across various regions and time periods. Your role is to assist in retrieving specific data as requested by the user and to provide additional analysis that highlights any interesting, unusual, or noteworthy aspects of the data, just as a human analyst would do.

When the user makes a request:

1. Accurately identify the relevant data source and retrieve the specific data they are asking for.
2. Perform a detailed analysis on the retrieved data to uncover any trends, anomalies, or key insights. Consider aspects such as:
  - Comparisons with previous periods or other regions
  - Significant changes or trends in the data

- Potential reasons behind the observed data patterns
- Any other insights that might be valuable for the user to know

Finally, present the data and your analysis in a clear, concise summary that the user can easily understand.

If the user's request is unclear or requires data from multiple sources, use your judgment to clarify the request and combine data sources as needed to provide a comprehensive analysis.

In this example, the objective is clearly described, as is what the agent should do if the end user asks it to do something outside of its prescribed scope.

### Multiple steps

Agents will execute multiple steps to meet their objectives. These individual steps are linked in chains, which define the steps the agent must take to meet the objective. This differentiates agents from the simple, direct prompting of an LLM. For example, asking an LLM to summarize a block of text cannot be considered an agent because it is a single step.

Take, for example, an agent that has been built to summarize financial reports. The multiple steps might be:

1. Call an API to download the desired report(s).
2. Locate and extract key figures from the report.
3. Look up historical values for these figures and compare them.
4. Extract key quotes from the report.
5. Generate a semitemplated summary that includes both extracted quotes, generated summary text, and comparison between historical and current figures.
6. Send the report to the recipient over the specified channel.

A step in the chain might involve a call to an LLM service with a templated prompt, or it could be the direct execution of a tool, such as an API call or a data formatting function. The steps are strung together in a chain, which may be sequential, looping, branching, or parallel chains. Throughout this process, multiple calls to the LLM service will occur without any user involvement.

## Autonomy

In the context of an agent, autonomy is its ability to plan, reason, and self-correct to achieve its objective with limited or no human intervention. For example, a minimal degree of autonomy for an analytics-generating agent may simply be deciding which Python package to use during the data analysis step. A more significant degree of autonomy may be choosing the tool that it will use to meet its objective from several made available to it (e.g., deciding if it should query historical data from a data warehouse or live data from a CRM to best respond to the user's request).

Less autonomy will mean that the agent is less flexible in the type of problem it can solve but more likely to give a good result in that narrower range. More autonomy will mean more flexibility, but more risk that the results will not be satisfactory. In the enterprise, agents are likely to be quite limited in their autonomy, with narrowly defined options available to them, especially during the early stages of their development and use. This may change over time as models and agent-building techniques evolve and improve.

## Tool use

A defining characteristic of an agent is its use of tools to accomplish its objectives. These tools are described in more detail in the following section.

### Thinking Like an LLM Mesh

As an object in an LLM Mesh, an agent expects some task as an input and is expected to provide a satisfactory result as an output. This broad definition reflects the breadth of what agents can be built to accomplish.

## Tools

In an LLM Mesh, a tool is any function or system that an agent is provided with to accomplish its task. As such, tools cover a very wide range of potential technologies. This breadth gives agents and agentic applications their incredible potential: they can automate and accelerate tasks, decisions, and operations that otherwise require manual work across the enterprise and its business systems.

The types of systems that can serve as tools in an LLM Mesh include but are not limited to:

- Internal data storage and retrieval systems, such as databases, data warehouses, and data lakes
- Enterprise software systems, such as CRM, human resources management systems (HRMSs), and enterprise resource planning (ERP) systems
- Advanced analytical assets, like predictive machine learning models (e.g., churn prediction, lead scoring, or fraud detection models)
- Programming and querying languages, like Python and SQL, along with specific packages or proprietary code
- External data APIs, such as financial data or weather services
- Other agents within the LLM Mesh: when one agent is used by another, it acts as a tool from the perspective of the calling agent, performing a specialized, encapsulated subtask without exposing its own internal logic, fulfilling the role of a tool in that specific interaction

For an agent to use a tool, it must understand what the tool is and how it works, so each tool must provide its own schema. This schema is what allows for some standardized interaction with the tool, despite the great diversity of tools that may exist in the LLM Mesh. For example, the schema for a tool that retrieves a customer's order status might be structured as follows:

- **Tool Name:** `getCustomerOrderStatus`
- **Description:** "Use this tool to get the current status of a customer's order. This tool requires a unique order ID."
- **Input Parameters:**
  - `order_id` (type: string, required: true): "The unique identifier for the customer's order."
- **Output:**
  - `status` (type: string): "The current status, which will be one of the following: 'Processing', 'Shipped', 'Delivered', or 'Canceled'."

This structured information allows an agent to understand what the tool does, what inputs it needs, and what kind of output to expect.

This schema-driven approach allows an agent to operate autonomously, which distinguishes it from a rigid, predefined workflow. Instead of following hardcoded steps, the agent consults tool schemas to dynamically plan its actions, giving it the flexibility to adapt to changing needs.

### Thinking Like an LLM Mesh

As an object within an LLM Mesh, a tool provides a schema, making itself available for use by an agent. The tool expects an input and provides an output as defined in that schema. Tools are very flexible and their schema is essential to their use by agents.

## Applications

In an LLM Mesh, an application is what makes an agent available to end users. The agent defines the logic that orchestrates the different objects from the LLM Mesh that are required to accomplish a specific purpose. The application is the interface and supporting functions that allow the end users to interact with the agent, to better understand the results provided by the agent, and to provide feedback to the developers. The application is also where certain services providing security, safety, and cost control are enforced.

There are several types of agentic applications, including:

- Chat interfaces where users interact with the agent iteratively
- Contextual assistants, either as desktop applications or browser extensions, that provide some additional functionality or assistance in the context where the user is working at that moment
- Backend or “headless” applications that run without direct end-user interaction

Agentic applications can have a wide range of interfaces and functionality. The abstraction and standardization within the LLM Mesh make it simpler for the developer to build the application in a way that clearly communicates to the end user how the agent underpinning the application is generating its results.

For example, in the case of an application that exposes an analytics-generating agent to end users, it will be easier for the end user to understand and trust the results if the application distinguishes between outputs that come from a query to a retrieval system or a tool versus outputs that are the result of the LLM's interpretation or suggestion. Furthermore, the end user will also be more likely to trust results if they can verify that the sources used and the query that the LLM generated are appropriate for the objective of the application.

Feedback mechanisms should also be built into the application to ensure that when an agent does not behave as expected, end users can flag this anomaly to the developers so that they can monitor the agent's performance and take corrective action if necessary.

In addition to direct user feedback, the LLM Mesh should provide automated monitoring services to proactively detect performance degradation or unexpected behavior. These capabilities, which are essential for maintaining the application, will be explored in detail in **Chapter 4**.

#### **Thinking Like an LLM Mesh**

Within an LLM Mesh, the application object includes the application itself, versioning for the deployed application, and logging of user interactions with the application.

The previous sections defined the building blocks of agentic applications. The challenge, however, is managing these diverse objects as they proliferate across an organization. To enable discovery, promote reuse, and enforce governance, an LLM Mesh requires a central system of record—a catalog—for all of these assets.

## **Cataloging LLM-Related Objects**

As an organization begins developing more agentic applications, the number of different objects it will need to use to build those applications will grow rapidly. This could become difficult to manage, with users hunting for different objects, re-creating existing objects, or using unapproved objects.

Overcoming these challenges starts by creating a central catalog for all of these objects. This catalog is a fundamental component of an LLM Mesh. The catalog should:

- Account for all LLM-related objects that are available for use in the enterprise
- Provide documentation that describes and provides instructions for using each object
- Track the version or other details about the ownership and development history of the object
- Assign a unique ID to each object to allow it to be referenced and tracked unambiguously

This information is stored in a structured format that allows human and machine discovery of the available objects. Having a central catalog of the objects provides various benefits for organizations as they begin building more agentic applications. Those benefits include:

#### *Standardization*

Only approved objects can be added and used after a vetting process.

#### *Governance and compliance*

You can maintain full transparency and traceability of which data is used with which LLM for which purposes, enabling business alignment and regulatory compliance.

#### *Security*

The catalog allows access controls to be defined and enforced, controlling which end users and automated systems have access to which objects.

#### *Composability*

Once registered, objects can be easily added to new applications where they are combined with other objects, accelerating the development process.

#### *Efficiency*

Less time is spent manually connecting different objects, accelerating application development.

Importantly, this catalog will be useful for both the end users and the LLM-powered agents that they will be building. The agents will also rely on the documentation to discover and use the objects, and the agents will be subject to the security model.

## Conclusion

In this chapter, we have learned about the various objects of an LLM Mesh, how they are abstracted, and how they can be integrated with one another. The final chapter of this guide will go into greater detail about how a specific agentic application can be built using an LLM Mesh. Before getting to that, however, the following chapters will describe the various federated services that an LLM Mesh must also provide to meet enterprise security, reliability, and cost requirements for the many agentic applications that will be built within it. [Chapter 3](#) will start with that most important of enterprise considerations: cost. How can the overall cost of an enterprise's LLM use be optimized? Read on to learn more.

# Quantifying and Optimizing the Cost of LLMs in the Enterprise

At the beginning of [Chapter 2](#), we established the importance of an organization's ability to develop its own agentic applications. These built-for-purpose applications will allow the organization to differentiate itself from its competitors, pulling ahead in the market. That said, developing and running them will come at a cost.

[Chapter 2](#) established how an LLM Mesh will simplify and standardize the development of these applications, thus reducing the cost of their development. However, agentic applications will also incur costs as they run. Minimizing costs while maintaining the required level of performance will allow an organization's budget to go further, supporting the deployment of these built-for-purpose agentic applications across more functions of the organization, thereby getting more return from the same investment in AI.

It is critical not to consider cost in a vacuum; rather, it must be considered in conjunction with performance (defined here as the quality, speed, and reliability of an application's output, a topic covered in detail in [Chapter 4](#)). If an organization were to drive down its spending on AI blindly (for example, by using smaller models or cheaper services), it could find that its agentic applications are no longer delivering adequate performance.

On the other hand, if an organization had a policy of always using the highest-performing model, it would likely end up overspending for a level of performance that is, in fact, not needed. The answer is to find the balance where cost is minimized for an adequate level of performance. The necessary level of performance, both in terms of speed and quality of output, will depend on the specific requirements of each application. Thus, the goal is neither minimal cost nor maximal performance but rather minimal cost while delivering the required performance.

A key capability of an LLM Mesh is providing the federated services required to measure and track cost and performance across the many different components used to build agentic applications. By using these federated services in an LLM Mesh, organizations can fully understand where their AI budget is going and enforce policies to get the most return on that investment.

We will tackle the topics of cost in performance in both this chapter and in [Chapter 4](#). In this chapter, we'll first understand the drivers for the cost of the different objects in an LLM Mesh, especially the different types of LLM services. Then, we will look at techniques for limiting costs. Finally, we'll consider the organizational practices needed to run a cost-efficient AI practice using an LLM Mesh, covering topics like cost reporting, budgeting, and rebilling. Throughout, we'll refer to the need to consider performance trade-offs when making decisions to reduce costs. Rest assured that we will bring much more precision to performance measurement and monitoring in [Chapter 4](#).

## Quantifying the Costs of Agentic Applications

In this section, we will understand how agentic applications generate costs. We'll start with a quick review of the objects of an LLM Mesh to determine which are cost-generating and which are not. From there, we'll then dive deep into the costs of the different types of LLM services, including those provided by model developers, CSPs, and those that you manage yourself. Finally, we'll compare two different agentic applications to see how choosing different LLM services impacts their costs.

# The Objects in an LLM Mesh That Drive Costs

Let's begin by breaking down the cost structures of the different objects in an LLM Mesh. Recall from **Chapter 2** the seven main objects of an LLM Mesh: LLMs, LLM services, retrieval services, prompts, agents, tools, and applications. Of these seven objects, some are cost generating and others are not.

## Non-cost-generating objects in an LLM Mesh

The following objects function as logical assets whose costs are tied to development and maintenance rather than direct, per-use execution:

### *Prompts*

Like code or documentation that an organization would develop and manage, prompts themselves have no direct cost. Similarly, while prompts do not have a direct, per-use cost, their development and maintenance represent a significant labor cost, especially as prompt libraries scale and require versioning.

### *Agents*

As logical objects, agents themselves do not generate direct costs. While the application that executes an agent's logic has hosting costs, the primary costs generated by an agent's operations are from the LLM services, retrieval services, and tools it utilizes.

### *Agentic applications*

Like the agent, the code of an application generates no direct costs. While the application code itself does not generate per-use costs, the infrastructure on which it is hosted and the services it calls both do. Furthermore, inefficient code can lead to higher computational costs.

## Cost-generating objects in an LLM Mesh

Conversely, the following objects generate direct costs through compute, hosting, or per-use API fees:

### *LLMs*

LLMs generate costs in two primary ways, depending on their deployment model. When self-hosted, a model incurs direct costs for storage and the compute infrastructure required to run it. When accessed via a managed API, these infrastructure and

licensing costs are bundled into the per-use service fee. Some proprietary models licensed for private deployment, like those from Cohere or AI21 Labs, also require a separate licensing fee in addition to hosting costs.

### *LLM services*

LLM services are the workhorses of agentic applications and their main cost drivers. Several different cost models for LLM services will be explored in more detail in the following sections.

### *Retrieval services*

Retrieval services generate costs per use, much like the LLM services that often underpin them.

### *Tools*

The cost profile for tools depends on their function. Many traditional tools, such as a standard data API, have fixed subscription fees or simple usage-based pricing that is already managed by existing enterprise solutions. However, an important exception exists for dynamic tools within agentic workflows. A tool that itself calls an LLM, runs a retrieval pipeline, or recursively calls other services will incur variable, usage-based costs. These costs are not fixed and must be tracked by the LLM Mesh's monitoring services.

While a holistic financial model must consider the total cost of ownership (TCO)—including infrastructure, engineering, and monitoring, which can be substantial—the cost services of an LLM Mesh focus specifically on the new and highly variable costs introduced by LLM and retrieval services. The costs of other components, such as traditional data querying or API services, are typically fixed or are already well-managed by existing enterprise solutions. The LLM Mesh is therefore adapted to the unique cost structures of the new AI assets in the enterprise IT landscape.

## **Additional costs in an LLM Mesh**

In addition to the costs of the objects combined to build agentic applications, organizations may incur the following costs:

- Licensing fees for any LLM Mesh services for analysis and control
- Licensing fees for a cataloging and documentation solution

- Licensing fees for other components of an LLM Mesh
- The costs of developing and maintaining any of the LLM Mesh components or infrastructure that they do not license from a software vendor; note that these costs, which may be significant, are out of the scope of the LLM Mesh and, thus, this guide<sup>1</sup>

Whether these components are licensed separately, as part of an all-in-one LLM Mesh solution, or developed internally will depend on the organization's LLM Mesh strategy, and each approach has different costs.

## Understanding the Costs of LLM Services

You must set up your LLM Mesh to capture the total costs of all agentic applications in a normalized way, allowing for their comparison and aggregation. By having a comprehensive and comparable way of looking at the costs of different LLM services, organizations can more easily build capabilities that combine models from different providers to strike the optimal balance between cost and performance across their entire fleet of agentic applications. Without the cost tracking and control services of an LLM Mesh, the costs of these different services are difficult to compare and control, making it challenging for organizations to optimize their spending.

As described in [Chapter 1](#), there are three main types of LLM services based on where the LLM is hosted:

### *Model developer–managed services*

These are the LLM services offered directly by model developers, such as OpenAI, Google Gemini,<sup>2</sup> Anthropic, Cohere, or Mistral. They are model as a service (MaaS) offerings, and costs are calculated on a per-token basis. Many developers will offer both on-demand pricing and reduced batch pricing.

---

1 This [article](#) and [calculator](#) provide a good basis for estimating the total cost of ownership for certain applications: Zilliz, “How to Calculate the Total Cost of Your RAG-Based Solutions,” Medium, February 23, 2025.

2 Google offers its Gemini model both directly to customers and through its Google Vertex AI product. The Google Vertex AI product offers the models of other developers as well. We will refer to the Gemini direct offering as “Google Gemini” and the offering that includes models from other developers as “Google Vertex AI” to avoid confusion.

### *Cloud service provider–managed services*

These are the LLM services offered by CSPs such as Amazon Web Services (AWS), Microsoft Azure, and Google Cloud Platform (GCP). The cloud platforms have given specific branding to their LLM services: Amazon Bedrock, Google Vertex AI, and Microsoft Azure AI Studio.<sup>3</sup> These services generally offer access to a number of curated LLMs that can be deployed on managed compute environments or as serverless MaaS instances.

### *Self-managed services*

Organizations may choose to manage their LLM services themselves. In this case, they may rent a server from a cloud provider or provision an on-premises server for themselves.

The following sections will explore the details of these different options and their pricing models, giving some examples of how those pricing models impact total costs. As you will see, there are diverse pricing models and options available, making the centralization and normalization of the cost service of an LLM Mesh an essential element for deployments that go beyond initial testing and experimentation.

## **The Cost Model of Developer-Hosted LLM Services**

The pricing model for the LLM services offered by the model developers mostly follows the same MaaS structure: a price per token input (i.e., the prompt) and a price per token output (i.e., the generated instructions). The price per output token is generally three to five times as expensive as per input token. **Table 3-1** summarizes the on-demand pricing of several developers' flagship models, showing the discrepancy between input and output tokens.

---

<sup>3</sup> Microsoft Azure's foundation model catalog is also accessible through its Azure Machine Learning Studio product.

Table 3-1. On-demand pricing of developers’ flagship models as of October 2024

Developer and model	Input tokens (per 1M tokens)	Output tokens (per 1M tokens)
OpenAI GPT-4o	\$2.50	\$10.00
Google Gemini 1.5 Pro (>128k context window)	\$2.50	\$10.00
Mistral Large 2	\$2.00	\$6.00
Anthropic Claude 3.5 Sonnet	\$3.00	\$15.00

These services will typically offer their latest models at a few different sizes, with larger models being more expensive and smaller models being cheaper. They may also offer specialized models, such as those for generating text embeddings, generating images, or generating code.

In addition to paying as you go, some model providers are beginning to offer options for reducing costs through two different approaches: offering lower pricing for asynchronous *batch processing* and reducing the input tokens needed through *context caching*.<sup>4</sup>

### Batch processing

While some applications, like chatbots, require immediate responses from LLM services, other applications do not. For example, using an LLM to classify the sentiment in a large dataset of historical customer reviews can be done asynchronously, with the request being sent and the response being returned sometime later. In the case of OpenAI, a 24-hour turnaround time for all batch jobs is guaranteed. Check with your provider to see if it provides reduced pricing for batch processing; it could be a valuable source of savings for the right kind of application.

When implementing an LLM Mesh, be sure to do so in such a way as to take advantage of batch processing where the LLM service offers it. And be sure that the LLM Mesh takes the price difference into account for cost tracking and analysis. This is done by configuring each service endpoint in the Mesh with its unique pricing model—such as

---

<sup>4</sup> As of October 3, 2024, OpenAI and Google Gemini offer both batch processing and context caching, and Anthropic offers context caching (called *prompt caching*). Cohere and Mistral do not currently offer these options.

different per-token rates for on-demand versus batch processing—and then applying the correct model to the usage logs for each call.

### **Context caching**

There are some cases where you will include repeated information in a prompt submitted to an LLM service. For example, a chatbot may include a lengthy system prompt at the beginning of each request to ensure consistent and appropriate behavior. The LLM service will retain a cache of these input tokens and can use that cache instead of the new tokens, resulting in reduced latency and cost. Some of the LLM services pass these savings on to their customers.

Google Gemini prices cached input tokens at a 75% discount for uncached tokens, while OpenAI offers a 50% discount for cached tokens. Both the OpenAI and Google Gemini services apply caching automatically—no change to your API calls is required to take advantage of it. That said, it is important to carefully read the developer’s documentation to understand under what conditions cached tokens will be used so that you can create the conditions to maximize these savings. For example, OpenAI only checks the first few tokens to determine if the cache should be used or not, meaning that you need to structure your prompts such that the repeating information is always at the very beginning.

The LLM services will report when cached tokens are used. An LLM Mesh should track this and tabulate costs accordingly.

## **The Cost Model of CSP-Managed LLM Services**

The major cloud providers have all introduced managed services that allow popular open source and proprietary models to be quickly deployed on infrastructure managed by the cloud provider. This can be a quick way for organizations with an existing relationship with one or more CSPs to test and begin building with LLMs from different developers.

In contrast to the offers from the model developers, which all follow a MaaS paradigm, the CSPs offer more pricing models and options. This introduces some additional complexity, but it also introduces more options and flexibility that allow an organization to optimize for cost and performance as well as data residency.

AWS, Azure, and GCP all offer two main pricing models: on demand and provisioned throughput.

### On-demand pricing

On-demand pricing is a MaaS offer that is billed on input and output tokens (like the MaaS offers from the model developers). In most cases, the models offered for on-demand pricing will come from several curated developers. For example, in addition to its own Gemini models, Google Vertex AI, an AI development platform for building and using generative AI, offers models from AI21 Labs, Anthropic, and Mistral. There is thus a trade-off between the convenience of the fully managed MaaS offer and the full choice of models available on model hubs, such as Hugging Face.

In contrast to the offers of the model providers described previously, only Amazon Bedrock offers batch pricing for its on-demand offer, **giving a 50% discount**. This pricing is only available for certain models in certain regions.

### Provisioned throughput

In contrast to their on-demand offers, the provisioned throughput offers from the CSPs reserve a defined amount of model capacity for a given amount of time. The primary benefit is guaranteeing resource availability for your applications. Thus, it is not primarily a cost-saving strategy but rather a way to guarantee the availability of required resources. This approach is well-suited for applications with predictable usage patterns, such as in B2B contexts where workloads are stable or have been tested over time. It is also a good fit for mission-critical tasks where consistent, low-latency performance is a requirement. Longer reservation periods result in a discount over the CSP's on-demand pricing.

In contrast to the per-token pricing of on-demand offers, provisioned throughput offers are less transparent. In a provisioned throughput offer, you reserve *model capacity*, which each CSP defines and names differently. On Azure, they are known as *provisioned throughput units*, and an **estimator tool is provided**, but not the underlying formula. On Amazon Bedrock, they are named *model units*, and you must contact your account manager to estimate the quantity that you would need.

Provisioned throughput offers are best adapted to applications that have already been tested and deployed to production and thus where the expected usage is relatively certain. If you have an application where you anticipate variable usage, the on-demand offers are probably the best bet.

CSPs generally do not provide as many models in their provisioned throughput offers as they do in their on-demand pricing offers. For example, Google only offers its Gemini models for provisioned throughput, and Microsoft Azure only offers models from OpenAI. Amazon Bedrock offers both its Titan family of models and **models from other developers**.

## The Cost Model of Self-Managed LLM Services

The previous two categories of managed services offer much convenience, but at the cost of control. While it is easy to begin using the managed services with a simple API call, you are limited in the models available and other aspects of the deployment. If more control over those details is required, then an organization may consider self-managing its own LLM services. This can be done by renting a server instance from a cloud provider or using an on-premises server, typically maintained by your organization's IT department.

Cost will be an important consideration when deciding whether to deploy self-managed LLM services. In contrast to the per-token or model unit costs of the managed services, the direct costs of the self-managed services will be the hourly cost of the server instances. The costs will thus be fixed, independent of the number of tokens you push through the service. A cost-conscious organization would thus seek to maximize the use of its service by maximizing token throughput, without overloading it and degrading its performance. A delicate balance must be struck, for sure!

The cost model for a self-managed LLM service will first depend on whether you are renting the instance from a cloud provider or using an on-premises server. In the case of an on-premises server, the cost model will depend on your organization's internal rebilling policies. In the case of instances rented from a cloud provider, two primary cost models exist: on-demand usage and savings plans that require a long-term commitment, typically one or three years. The following sections explore these in more detail.

The pricing of the server instances available to rent from the major cloud providers will depend on their combination of GPU memory, compute, and storage resources. GPU memory is generally the constraining resource that you need to take into account, depending on the size of the LLM that you would like to host. This is because all of the model's weights need to be loaded into the GPU memory. Very large models with tens or hundreds of billions of parameters will exceed the memory of even the largest single GPUs and must therefore be deployed to multi-GPU instances. To choose the right server, you will need to estimate the throughput required for the anticipated use of your application(s).

You will need to ensure that the server instance has all of the necessary drivers and software to serve the models. The cloud providers typically offer base images for these instances that provide the necessary drivers, but be sure to check the documentation from the model developer to make sure that you are not missing anything.

### **Costs of self-managed, on-demand cloud server instances**

On-demand servers are usually priced per hour, and you begin paying the moment they're activated, regardless of the number of tokens they're processing. In this way, it is a fixed cost. It is important to thus spin down any instances that are not in use.

Even though the cost of the instance is expressed per hour, most CSPs offer finer granularity for billing, often down to the second. So if you use the instance for 8 hours, 12 minutes, and 23 seconds, you will only pay for those 29,543 seconds.

That said, be wary of attempting to overoptimize your instances by spinning them up and down too frequently, as there will be a cost in terms of latency. It will usually take a few minutes for an instance and the connected services to come fully back online and to be available for use. Thus, "on demand" only refers to the lack of long-term billing commitment and should not be misinterpreted as meaning "available for immediate use." On-demand instances are most appropriate for predictable workloads that will remain constant for a period of time.

### **Costs of self-managed, long-term cloud server instances**

If you expect the workload of your application to remain constant over a year or more, long-term commitments can be very

cost-efficient. By agreeing to pay for a predetermined volume of server usage, the CSPs will usually offer a discount over the equivalent on-demand price. The discount rate will depend on the instance type, the region, and whether you pay upfront. Discounts can be as large as 70% but may be smaller in practice, particularly for in-demand instance types like those for LLM workloads.

Note that the commitment implies that the instance is paid for the entire duration, meaning that you cannot spin it down during periods of no use. Thus, a three-year commitment implies multiplying the hourly cost by the number of hours in three years (26,280, assuming there is no leap year).

### **Costs of self-managed, on-premises server instances**

Some organizations, though fewer and fewer all the time, prefer—or are required—to manage their server infrastructure. If your organization operates this way, you are undoubtedly aware of it already, and you are undoubtedly aware of the process of gaining access to these resources. In this case, a guide like this will not be able to provide you with detailed information, as the costs and how they are rebilled will depend on your organization's practices.

In all cases, an LLM Mesh should allow you to move easily among the services available to your organization and must capture the costs of all different types of services. As the pricing of many hosted services is not available programmatically (i.e., they cannot be looked up via API), an administrator of your LLM Mesh should track them and update them carefully to ensure that the cost data is captured accurately.

## **Comparing the Costs of Applications**

To make these differences more real, let's imagine two different applications in a large enterprise. These applications will have very different usage scenarios to show the consequences of using different applications with different services.

### **Application 1: Company-wide knowledge assistant**

This first application is a knowledge assistant application that uses a RAG pipeline to surface relevant information to the user from an internal document base. We imagine this application will be deployed in a large, global enterprise. As a result, it sees relatively

constant usage across time zones and throughout the year. As a RAG application, a large volume of text extracted from the source documents is passed to the LLM service in the prompt, resulting in a large volume of input tokens. Let's build out a low-volume and a high-volume scenario to estimate the total token count for the LLM service over a year; those estimates are shown in [Table 3-2](#).

*Table 3-2. Usage hypotheses for company-wide knowledge assistant application*

	Low-usage scenario	High-usage scenario
Input tokens per session	5,000	15,000
Output tokens per session	500	5,000
Sessions per user, per day	5	20
Users	1,000	1,000
Working days per year	250	250
Total input tokens per year	6,250,000,000	75,000,000,000
Total output tokens per year	625,000,000	25,000,000,000

Now, let's compare the costs of running this application against two different LLM services. The first is an on-demand service from a model provider. We'll use OpenAI in this example. The second is a self-managed cloud server from AWS. This type of application does not need the advanced reasoning capabilities of the largest, most sophisticated models. As such, we'll choose the OpenAI GPT-4o mini as the MaaS option. This scenario is detailed in [Table 3-3](#).

*Table 3-3. Pricing scenario for OpenAI GPT-4o mini on demand as of October 2024*

	Low-usage scenario	High-usage scenario
Total input tokens per year	6,250,000,000	75,000,000,000
Total output tokens per year	625,000,000	25,000,000,000
Input price (per 1M tokens)	\$0.15	\$0.15
Output price (per 1M tokens)	\$0.60	\$0.60
Total input cost	\$937.50	\$11,250.00
Total output cost	\$375.00	\$15,000.00
<b>Total cost</b>	<b>\$1,312.50</b>	<b>\$26,250.00</b>

For the self-managed option, we will select, just for the sake of the example, the Llama 3.2 11B model and will run it on an AWS g5.2xlarge EC2 instance. This instance should be sufficient for the hypothetical usage, but it is important to model the required resources appropriately for your expected use, as different instance configurations have very different prices. In this scenario, we will compare the costs for no commitment (on demand), a one-year commitment, and a three-year commitment. As these costs are fixed regardless of usage, they would be identical for the low- and high-usage scenarios; hence, we show cost in columns per commitment period (see [Table 3-4](#)).

*Table 3-4. Cost for AWS EC2 g5.2xlarge (US East, Ohio) as of October 2024*

	No commitment (on demand)	1-year commitment	3-year commitment
Cost per hour	\$1.212	\$0.95445	\$0.65448
Cost per year	\$10,617.12	\$8,360.98	\$5,733.24

In most cases, self-managing the AWS EC2 instance would be less costly, with the exception of the lowest usage scenario. In the high-usage scenario, making a three-year commitment and self-managing the EC2 instance can deliver 78% savings. This shows the importance of estimating and monitoring your usage of LLM services and making the correct choice per application. Once again, these estimates exclude the significant overhead costs of self-management, such as engineering, monitoring, and security.

### Application 2: Corporate strategy sparring partner

Let's now imagine a very different application. This application is built on top of an agent that is designed to support the strategic planning activities of the senior leadership and corporate strategy teams. It is used far less per year, by far fewer people. That said, the usage is very intensive, and it requires a very large, very capable model. Let's go through the same exercise, starting with the usage hypotheses (see [Table 3-5](#)).

Table 3-5. Usage hypotheses for corporate strategy sparring partner application

	Low-usage scenario	High-usage scenario
Input tokens per session	5,000	50,000
Output tokens per session	5,000	25,000
Sessions per user, per day	10	50
Users	20	20
Working days per year	50	50
Total input tokens per year	50,000,000	2,500,000,000
Total output tokens per year	50,000,000	1,250,000,000

As with the previous application, let’s now consider the costs of running this application. To keep the comparisons similar, we’ll use OpenAI again as the MaaS provider, but we’ll select its flagship model, GPT-4o, shown in Table 3-6.

Table 3-6. Pricing scenario for OpenAI GPT-4o on demand as of October 2024

	Low-usage scenario	High-usage scenario
Total input tokens per year	50,000,000	2,500,000,000
Total output tokens per year	50,000,000	1,250,000,000
Input price (per 1M tokens)	\$2.50	\$2.50
Output price (per 1M tokens)	\$10.00	\$10.00
Total input cost	\$125.00	\$6,250.00
Total output cost	\$500.00	\$12,500.00
<b>Total cost</b>	<b>\$625.00</b>	<b>\$18,750.00</b>

Imagine now that you want to self-manage this application. You want to use a large and highly capable model, so you choose Llama 3.2 90B. Given the size of this application, you will need to use an EC2 instance with multiple GPUs, in this case, a g5.12xlarge. That choice results in the pricing shown in Table 3-7.

Table 3-7. Cost for AWS EC2 g5.12xlarge (US East, Ohio) as of October 2024

	On demand	1-year commitment	3-year commitment
Cost per hour	\$5.672	\$4.4667	\$3.06288
Cost per year	\$49,686.72	\$39,128.23	\$26,830.83

For an application that is only used for part of the year and requires a high-performing model, the most cost-efficient approach will be to use the MaaS offering from the model provider, even in the highest usage estimate. Paying for a self-managed instance of sufficient size for the model required would be wasteful.

We've intentionally constructed the high- and low-usage estimates to show the high potential variability between applications. For example, applications with an agent can cycle through many prompts and responses during CoT reasoning. This can result in a large volume of both input and output tokens that the end user never sees but which are necessary to the proper functioning of the application. Monitoring this spending is essential to make sure that the budget is used wisely and that costs are minimized.

## Techniques for Limiting Costs

In the previous sections, we reviewed how to measure and monitor costs in an LLM Mesh. Now, how can you limit costs, and how should you approach a new project in a cost-conscious way?

When building out a new agentic application, a good rule of thumb is to start by defining the level of performance that you need for a given application without much consideration for cost. This is because it is important to confirm if the application will work at all and not be left wondering if it would have worked if a more powerful and more costly model had been used instead. Thus, you should start with a large, high-performing LLM that will allow you to build out your application quickly. Then, once you have achieved the level of performance that you need, you can start optimizing your application, testing lower-cost models and techniques to maintain the level of quality and speed that your application requires while minimizing the cost. The following sections review some of these techniques.

## LLM and LLM Service Selection

The simplest method of reducing cost is to switch to a cheaper LLM service that provides the required level of speed and performance. There are two ways to think about this: model upgrade and model substitution.

## Model upgrade

Upgrading a model is when you move from one model to its successor: for example, upgrading from OpenAI's GPT-4 to GPT-4o. The good news is that, given the rapid rate of development in the space today, the new generation of a given model is often both better performing *and* cheaper than its predecessor. Such is the magic of living in a time of rapid technological progress, and while this trend may not continue, you should take full advantage of it.

An LLM Mesh should be configured so as to allow administrators to monitor for these upgrade opportunities and then facilitate the rapid testing to confirm that the new model does, in fact, provide improved application performance. You may ask, how could it not? While the new model will almost certainly outperform the previous model on the benchmark metrics, that does not necessarily mean that it will automatically improve the overall performance of the application. For example, your application may depend on very carefully crafted prompts that exploit some quirk of the outgoing model. The new model may not behave in precisely the same way, leading to degraded performance. Once you adapt the prompts to the new model, you can likely achieve equivalent, if not improved, performance at equivalent or lower cost.

Upgrading a model may seem like a trivial task. However, the many steps in the previous paragraphs demonstrate several of the benefits of using an LLM Mesh. First, with an LLM Mesh, developers can quickly introduce the new model—no application logic needs to be changed. Second, performance measurement in the LLM Mesh allows the developer to quickly measure the performance with the new model and to identify quickly which prompts or other objects are now seeing degraded performance, allowing corrective measures to be taken. Finally, cost tracking would show the difference in the cost between the application with the old model and the application now with the new model. By reducing the effort required for each of these tasks, an LLM Mesh makes it easier and faster to improve agentic applications, integrating improved models as quickly as they become available.

## Model substitution

In addition to upgrading to a new version of the same model, a viable option may be to consider a different model or family of models. Here, the decision will not be triggered by the release of a

new model. Instead, a team will seek opportunities to reduce the cost of its application, testing lower-cost options to see if they can be used while maintaining the quality required for the application.

### **Service substitution**

Finally, in certain situations, consider migrating from one LLM service to another, even if the underlying model is the same. For example, imagine an internal knowledge management chatbot that sees regular, predictable usage throughout the year. Your organization anticipates that it will continue to use this application for several years without much modification. If it runs on a CSP-managed instance with its on-demand option running an open weights model (Llama 7B, for example), you could consider migrating to a self-managed service running on a server instance with a long-term commitment to reduce your costs. Note that it will now be incumbent upon your organization to manage this service. Such a migration is a nontrivial engineering task that requires carefully matching the model's version, tokenizer, and configuration to ensure consistent performance.

## **Prompt and Inference Optimization**

Whether using a managed service or self-managing your LLMs, optimizing your prompts to use fewer input and output tokens while delivering the desired results is one of the most direct ways to reduce costs. Well-established manual techniques (popularly known as *prompt engineering*) as well as emerging programmatic techniques exist for optimizing prompts. Additionally, an LLM Mesh can help optimize inference through caching. The following sections describe these techniques.

### **Manual prompt engineering**

The principle of manual prompt optimization is to craft better, often shorter, prompts. Many guides are available online, often providing helpful advice. An overview of prompt engineering techniques is beyond the scope of this guide, but in general, being direct and concise in the instructions and clearly specifying the expected output will help minimize the volume of input and output tokens per request.

Remember, in an LLM Mesh, prompts are objects that are developed, tested, registered, and reused. Thus, optimized prompts that deliver good results at low cost can be shared and reused across multiple applications, reducing the time needed for application developers to test new prompts. In an LLM Mesh, a prompt must be linked to a specific model, as not all prompts perform equally across all models.

## Prompt compression

If LLMs are good at generating text, and a prompt is just text, can we ask an LLM to improve a prompt? The answer is yes, and there is much active research happening at the frontier of programmatic prompt optimization. These methods generally use an LLM to analyze a prompt, to understand which tokens in the prompt are actually driving the desired result from the model, and then compress the prompt down to these essential tokens. The result is often a prompt that is unintelligible to a human reader but delivers the desired results with far fewer tokens. These approaches are particularly useful in prompts that provide a lot of context to the model, such as those in RAG pipelines, where a large volume of text is sent to the model.

Two popular approaches to prompt compression are LLMingua,<sup>5</sup> a research project developed by Microsoft, and **Selective Context**, an open source project developed by academic researchers. These approaches can allow up to a 20× compression rate (i.e., reducing input tokens by 95%) while maintaining the required performance.

Implementing prompt compression within an agentic application, particularly one with long internal prompts—such as a RAG pipeline or an agent using a CoT reasoning technique—can be a powerful strategy for reducing the input tokens to the LLM service, thereby reducing costs.

By offering prompt compression as an option, an LLM Mesh ensures that application developers don't spend time repeatedly implementing such techniques.

---

<sup>5</sup> *Microsoft Research* (blog), “LLMLingua: Effectively Deliver Information to LLMs via Prompt Compression,” 2023, [https://oreil.ly/Gv\\_Ns](https://oreil.ly/Gv_Ns).

## Caching

As described in this chapter, some LLM services offer reduced pricing for prompts where context caching can be applied. An LLM Mesh can be configured to detect opportunities where a prompt can be adapted to trigger context caching when using a service that proposes this option. Usually, this means making sure that prompts with similar content start with identical strings of text.

In addition to ensuring that context caching is triggered when possible, an LLM Mesh can provide its own caching capabilities. For example, an LLM Mesh can detect when a prompt being sent to a particular LLM service is identical to a prompt sent recently. Rather than send the new prompt, it can simply supply the previously generated response, avoiding the unnecessary regeneration of the response and reducing both latency and cost.

## LLM Modification

In the previous sections, we have looked at ways to reduce costs without changing the LLM itself. Those techniques can thus be used on proprietary models or with LLM services that don't allow for model modification. A more advanced approach is to modify the model itself so that it can provide the required results at a lower cost. The following sections review those approaches.

### Model quantization

In model quantization, the precision of the model weights is reduced by converting them typically from 32-bit floating-point numbers to 16-bit floating-point or 8-bit integers. In layperson's terms, the more precise values are rounded to a less precise value with fewer significant digits. This process reduces the model's size on disk and simplifies the calculations required for inference. The result is lower latency and reduced computational cost, making quantization an effective technique for deploying models of any size onto smaller devices. This performance gain, however, represents a trade-off with model accuracy. As such, the results of a quantized model must be tested to ensure that they still meet the necessary quality requirements.

An LLM Mesh can offer model quantization as an option for any LLM objects. Note that this is only relevant for open-weight models running in a self-managed LLM service.

## Model pruning

While model quantization reduces computational intensity by lowering the precision of weights, model pruning does so by reducing the number of weights in the model. This approach is particularly effective when the goal is inference efficiency (lower latency and cost per token), when working with a domain-specific model instead of a generalized one, or when a model is overparameterized for its task.

Pruning is also a valuable strategy to use before fine-tuning a model. Emerging techniques, such as LLM-Pruner,<sup>6</sup> developed by a team of researchers from the National University of Singapore, show promising results. Although some initial performance degradation can occur after pruning, this is often temporary and can be recovered by fine-tuning the model.

As with quantization, an LLM Mesh should offer model pruning as an option for any applicable LLM objects. Note that this is only relevant for open-weight models running in a self-managed LLM service.

## Fine-tuning

A much-discussed approach to improving the cost efficiency of an LLM is to fine-tune it to a particular context. From a cost perspective, however, the fine-tuning process will entail its own costs. But it can be less costly to use a fine-tuned LLM in the long run if, in order to get the results that you need, you find yourself having to provide many examples to the LLM in your prompts. By fine-tuning, the LLM can permanently “learn” to give the results that you want, reducing the cost per use from that point.

The fine-tuning process itself can range from updating the full model’s weights to more modern, parameter-efficient methods like low-rank adaptation (LoRA) that modify only a small subset of the model’s parameters on a specialized dataset. In the example of marketing copy, this would be reference texts that exhibit the desired style. If the fine-tuning is successful, the model will be able to correctly mimic the style without requiring multiple examples in the prompt. While there is a fixed cost to the fine-tuning process,

---

<sup>6</sup> Xinyin Ma et al., “LLM-Pruner: On the Structural Pruning of Large Language Models,” preprint, arXiv, September 28, 2023, [https://oreil.ly/t07t\\_](https://oreil.ly/t07t_).

this may be recouped in the reduced run costs using this model for this application.

When deploying an LLM Mesh, be sure to provide methods that allow for fine-tuning as an option for any open-weight LLM object. This way, the developers using your LLM Mesh can use the approach without implementing the process themselves.

## Cost-Efficient AI Operations in the Enterprise

In the previous sections, we learned about what drives the cost of agentic applications, how performance can be measured and balanced against cost, and many different techniques that can be used to reduce cost while maintaining the required performance.

However, running a cost-efficient AI practice requires more than knowing about or accessing the best cost-reduction techniques. It requires organizational policies and practices that ensure that those techniques are fully applied. This section is about these organizational aspects.

### Tracking and Reporting Costs and Performance

A major advantage of an LLM Mesh is that it allows you to track the cost and performance of all of your agentic applications in a single place. The centralization of this tracking is an essential part of a well-governed generative AI practice. Without the standardization that an LLM Mesh provides, however, you would need to manually aggregate this information from the many different applications being developed across the organization, each built in a heterogeneous way. Doing it this way would be a major barrier to obtaining a single view of all cost and performance data.

With regard to cost data in particular, be sure that the LLM Mesh that you deploy tracks the costs in a fine-grained manner, allowing for the costs to be aggregated across multiple dimensions, such as:

- Individual users
- Teams
- Departments
- Projects
- Functions

- Business priorities
- Usage type (experiments, development, production)
- Geography and region
- LLM provider
- LLM type
- LLM (version specific)
- Hosting architecture (model provider, CSP[s], self-managed)

By associating costs with all of these dimensions, you can generate the reports required by leadership to understand where the AI budget is being spent and to what benefit. It will also allow you to quickly identify the root causes of any anomalies. For example, if the costs of several different projects spike and they all use the same model, it could be that a change in the model has resulted in degraded performance for certain shared prompts, which would then need to be corrected.

Reports across these dimensions can be used to support a culture of transparency and accountability for all AI costs. Specifically, these reports can be shared with:

- The developers of the application, to help them to understand the ultimate cost of the applications that they are developing and to help make them aware of the consequences of their design choices
- Management and budget owners, supporting a culture of transparency and accountability for all AI costs

## Setting and Enforcing Budgets

Based on the tracking and reporting capabilities described in the previous section, an LLM Mesh can also allow you to set and enforce budgets.

From a technical perspective, setting a budget is simple: you set a value that should not be exceeded for a given cost dimension or combination of cost dimensions. From that number, an LLM Mesh can allow you to enforce that budget in several ways.

Here are some examples of how that budget can be enforced in increasing order of severity:

#### *Warnings*

The LLM Mesh can alert the budget owner that they will soon reach or have reached their budget. This ensures that they are aware of the issue and can take appropriate action. Best practices have alerts set at 50% and 90% of the budget.

#### *Throttling*

If a budget has been exceeded, the LLM Mesh can throttle cost-incurring traffic to LLM services. This will help limit the exceedance without entirely interrupting service. That said, it will degrade service and thus may not be appropriate for mission-critical applications.

#### *Blocking*

Beyond throttling, an LLM Mesh can also be configured to block certain LLM services if the budget has been exceeded. This could be useful, for example, in the case of an experiment that has gone wrong: The developer may not be aware of the costs they are generating, and an automatic block can prevent a costly and embarrassing overrun. On the other hand, it would be inappropriate to block a customer-facing application. Blocking should be context aware, and graceful degradation must be in place.

## **Rebiling Policies**

Organizations may pursue a policy of rebilling the cost of running an agentic application to the business unit that benefits from its use. An LLM Mesh can support this practice in the following ways:

#### *Transparency*

The business units that bear the cost can clearly understand how the application works and which parts drive the cost.

#### *Reassurance*

The business units can be informed about the cost-limiting techniques that have been applied, reassuring them that the application has been implemented in the most cost-efficient way possible.

### *Trust*

The business unit can trust that the costs are captured accurately.

### *Fairness*

The business unit can be reassured that other business units are also bearing the costs of the applications that benefit them and that the costs are being calculated consistently.

## **Defining and Enforcing Cost-Saving Policies**

As your organization builds more and more agentic applications using an LLM Mesh, you should define cost-reducing policies. These policies can take a variety of forms:

### *Recommendations and best practices*

As part of your training guidance for using the LLM Mesh, you can provide recommendations and best practices for all the cost-saving techniques that can be applied.

### *Periodic application reviews*

Once deployed and running, the LLM Mesh enables you to review the cost profile of the applications using the centralized cost tracking and reporting to supplement periodic manual checks and audits.

### *Mandatory approvals*

Within an LLM Mesh, you can enforce a mandatory cost review and approval process to ensure that the best practices have been applied before deploying an application to production.

### *Automated detection of cost-saving opportunities*

As an LLM Mesh is aware of all applications and tracks their costs, it can automatically detect opportunities to reduce the costs of the applications. The LLM Mesh could then raise an alert to trigger a review and decision about whether to apply that technique.

### *Automatic application of cost-saving techniques*

Certain techniques can be applied automatically in the background to all applications. For example, request caching could be used to automatically store and retrieve responses for identical prompts, avoiding the cost of reprocessing a duplicate request.

# Conclusion

There is no guarantee of which approach will result in the best combination of cost, speed, and quality when developing and deploying agentic applications. As such, the best that any organization can do is to test the many different approaches to optimize its applications, to learn from this experimentation, and then seek to generalize the best practices as policies.

An LLM Mesh supports this in several ways:

- An LLM Mesh makes the different cost-saving *techniques* easily available to all application developers. This ensures that an individual developer does not waste time implementing these methods themselves.
- An LLM Mesh provides *visibility* into all costs so that they can be measured in development and monitored during deployment, facilitating reporting and budgeting.
- An LLM Mesh allows cost-saving policies to be *enforced*, ensuring that best practices are respected.

As such, an LLM Mesh allows cost-efficiency to become a core strength of an AI practice. This means that an organization can successfully develop more agentic applications in more business domains for the same budget. This cost efficiency is essential to maximize the value that your organization derives from generative AI.

But life is not without trade-offs. These cost-reduction techniques may degrade the performance of your agentic applications. In the next chapter, we will learn about how to measure and monitor performance so that you can maintain the required level while keeping costs at a minimum.

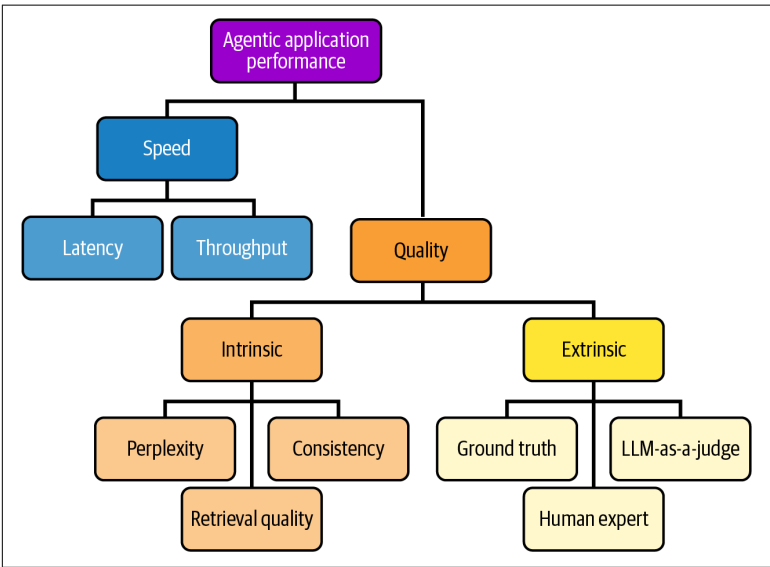
# Measuring and Monitoring the Performance of Agentic Applications

Now that we understand the cost models for various LLM services, let's look at how the performance of LLMs—and the agentic applications built on top of them—can be measured. Remember, the objective of all of this is to be able to define the required level of performance for a given application and then find the lowest-cost way to deliver that level of performance.

But first, what do we mean by *performance*? Indeed, there are two very different notions here:

- The *quality* of the generated response. In other words, how well the generated response corresponds to the requirements of the application. This is the main focus of this chapter, as measuring the quality of LLM responses is a novel field where an LLM Mesh can provide significant value to an application developer.
- The *speed* of the service, in terms of how quickly a response is generated. We'll touch on this briefly at the end of this chapter as this monitoring is similar to established DevOps practices for monitoring the speed and responsiveness of API services. An LLM Mesh does not need to do more than capture these metrics.

Within the two dimensions of quality and speed, there are several subdimensions as well. **Figure 4-1** illustrates these dimensions and subdimensions as a tree diagram.



*Figure 4-1. Dimensions and subdimensions of measuring the performance of agentic applications*

The sections that follow will explore these dimensions in more detail. But let's start by describing how measuring performance fits into an LLM Mesh.

### **Are LLM Benchmarks the Same as Agentic Application Performance?**

When we talk about the performance of an agentic application, it is important to note that we are *not* talking about LLM benchmarks. LLM benchmarks are tests to which developers submit their models to compare their inherent performance against one another. However, the benchmark landscape is fragmented, with no single set of standards adopted by all developers. Some common benchmarks include Massive Multitask Language Understanding (MMLU), Instruction-Following Eval (IFEval), and Graduate-Level Google-Proof Q&A (GPQA), which generally involve submitting a standard list of tasks to a model and comparing the results.

These benchmarks are useful for comparing one model to another and thus can help choose which model you want to start building with. However, they do not measure the quality of the output of an agentic application for your specific needs, nor do they measure the responsiveness of your LLM service. Metrics that measure the quality of output and responsiveness of your service are critical for managing a fleet of agentic applications and are thus the focus of this chapter.

An LLM Mesh-based approach to evaluation improves upon the limitations of the alternatives: fully distributed evaluation and fully centralized monitoring. While each approach has its advantages and disadvantages, implementing evaluation within an LLM Mesh will provide many long-term benefits. [Table 4-1](#) compares these approaches.

*Table 4-1. Comparing architecture paradigms for evaluating the performance of agentic applications across the enterprise*

	Fully distributed evaluations	Centralized monitoring via an application	Evaluation in an LLM Mesh
Description	An evaluation system is built for each agent and application.	A centralized application captures logs to provide centralized monitoring.	Evaluation is provided as a service for all LLM services and agents.
Advantages	They're easy to kick-start and customize.	It matches the existing paradigm for application performance monitoring.	Evaluation can be dynamically leveraged by applications, agents, and multiagent systems. Distributed evaluation scales with distributed development. It maintains a common evaluation framework company-wide. Reusability of evaluation components promotes standards and efficiency.
Disadvantages	Central IT lacks the ability to set up diagnostic and remediation procedures for agentic applications. There is a lack of enterprise global quality standards.	Evaluation results cannot easily drive adaption and behaviors. Monitoring probes monolithic applications rather than reusable components.	There's a new application paradigm for the enterprise.

When implementing an LLM Mesh, you should:

- Provide a shared service for performance measurement and monitoring
- Allow a developer to use that service flexibly at different levels in their application

To begin, it is important to provide performance measuring and monitoring tools as a shared service for two reasons. The first is efficiency: you don't want your application developers to develop and redevelop similar common capabilities across many different applications. The second is consistency: you want to be able to compare performance across different LLMs and agentic applications. If every developer is implementing their own performance measurement metrics in their own way, it will be difficult, if not impossible, to gather consistent performance metrics across your growing fleet of agentic applications.

Then, it is important to allow your application developers the flexibility to apply the performance measurement and monitoring service of the LLM Mesh at different levels of their applications throughout the many calls to an LLM that an agentic application will make. When implementing an LLM Mesh, you should ensure that the evaluation can be made at the level of the call to the LLM service and not higher up in the application stack. By bundling the LLM service call and the evaluation call, it means that developers need to call only one API and that performance evaluation can be deployed consistently across all agentic applications.

The concepts of performance and quality in an LLM Mesh are similar to that of data quality in a data mesh. Quality metrics for agentic applications should be understandable and readily available because they are a key part of the “contract” that an agentic application has with the entity using it. In the same way that a data mesh is designed to establish and enforce such a quality contract so that end users are sure that they can trust the data that they are using, an LLM Mesh does the same by establishing a performance contract for agentic applications.

As we will see, reliably assessing the quality of an agentic application will require a combination of different techniques. An LLM Mesh should provide these different techniques as shared services that developers can experiment with and use freely without having to implement them themselves for each and every application they are developing. As with the methods for controlling cost discussed in [Chapter 3](#), providing pre-implemented performance monitoring methods to all developers of agentic applications in your organization as a shared service is a critical part of implementing an LLM Mesh. This way, they continue their focus on their applications and not on common components like performance monitoring.

## Measuring and Monitoring the Quality of Generated Text

There are three phases in the lifecycle of an agentic application where the quality of the generated text should be monitored:

### *Predevelopment*

Defining the quality metrics and their required levels

### *Development*

Measuring the quality of the results and iteratively improving the application to improve the quality metric

### *Deployment*

Monitoring the quality metric for changes and taking necessary measures to fix issues that arise

In the past, if you were developing a predictive ML model, proceeding through these steps would be rather straightforward. You would choose a metric depending on the business requirements of the application (e.g., in certain applications, you may be more sensitive to false positives than false negatives or vice versa) and set a minimum threshold for that given metric. You would then train the model to deliver the required performance. Finally, you would deploy the model and set up a monitoring capability to ensure that it continues to deliver the required level of performance in the face of real-world data. These processes are now well-established with standard best practices to follow.

Measuring the quality of agentic applications is very different for two reasons:

*Model outputs are nondeterministic*

With a traditional ML model, the same inputs always give the same outputs. In this sense, it is deterministic. With an LLM, given the same input and prompt, the agent can generate different outputs.

*The models are used in open-ended contexts*

In a given application, the same model might be called on to select a tool, generate code, evaluate a response, and generate a text response for the user, requiring different quality metrics and requirements for each interaction.

Both of these characteristics are features and not bugs of LLMs. The fact that they are nondeterministic means that they can mimic creativity and come up with novel responses. The fact that they can be used in open-ended contexts means that they can be flexibly used to solve many different problems. These are their strengths, but it means that measuring their quality is much more difficult than measuring the quality of predictive models, where statistical methods alone can be used to indicate whether a model is performing well or not.

Evaluating the quality of LLM outputs is done in two ways:

- Measuring how well the model is performing independent of any particular task—known as *intrinsic quality*
- Measuring how well the model is satisfying the task at hand—known as *extrinsic quality*

Both approaches are useful in measuring and monitoring the quality of LLM outputs. Measures of intrinsic quality are useful to identify problems in model performance early on in a relatively low-cost way. At the same time, measures of extrinsic quality are necessary to ensure that the output of the model is actually solving the task at hand in a way that a human expert would judge as appropriate. When implementing an LLM Mesh, both intrinsic and extrinsic evaluation techniques are required.

## Intrinsic Quality Evaluation

Intrinsic quality measures assess how well the LLM is performing, independent of the task at hand. These are different from the LLM benchmarks mentioned in the note box at the beginning of the chapter. Those benchmarks are actually extrinsic measures of how well an LLM performs at a standardized task. Instead, the intrinsic measures look at several aspects of the inherent performance of the model. Those include:

### *Perplexity*

How confident the model is in predicting each token

### *Consistency*

How similar model outputs are when provided with the same input

### *Retrieval quality*

How effectively an underlying retrieval system is finding relevant documents

Let's now look at each of these three factors.

## Perplexity

Perplexity measures how “surprised” a model is by a sample of text, which reflects how well the model's predictions match that text. A low perplexity score means the model finds the text more predictable, assigning a higher probability to that sequence of words. However, this predictability does not equal factual accuracy. A model may assign a high probability to an incorrect statement, producing a plausible-sounding but wrong answer.

While perplexity does not guarantee accuracy, high perplexity is a useful indicator that the model is entering what is, for it, unfamiliar territory. In an agentic application, this could be caused by an out-of-distribution input, a mismatched context, or an imprecise prompt. Although high perplexity can correlate with hallucinations, it does not mean the output is necessarily wrong; it could simply be encountering a new topic. However, a practical challenge is that many commercial LLM APIs do not expose the per-token likelihoods required to calculate perplexity, limiting its use as a real-time metric.

For example, if high perplexity is measured at a step where the agent chooses from a list of available tools, it indicates that the model is not sure of its choice and may have chosen a different tool in another, similar case. By monitoring perplexity, you can detect this situation and take steps to fix the problem. For example, you could improve the prompt by providing more detail to the model about the task at hand, or you could change the schemas of the tools so that the model better knows which one to choose.

Another important aspect of measuring perplexity is that it can provide input to the model itself about the best next step to take, enabling the kind of agentic feedback loop described in [Chapter 2](#). For example, by feeding an indication of perplexity back to the model, the agent can be instructed to escalate the task for human review or stop an automatic process. The specific threshold for escalation will depend on the specific use case and will require testing to identify. These are important performance and safety measures that are only possible if perplexity is measured consistently.

## Consistency

Even if a model has a low perplexity score for its response, an enterprise application still requires consistent performance. It is, therefore, important to measure the model's consistency by providing it with the same input multiple times and measuring the similarity of the outputs. While traditional NLP techniques like cosine similarity can be a starting point, they can be brittle for evaluating the nuanced, open-ended responses from LLMs. More advanced semantic similarity metrics are often required to capture true consistency.

By continually monitoring the consistency of a model within an LLM Mesh, you can detect early on any problems that may affect the quality of any applications that are built on top of the model.

## Retrieval quality

While perplexity and consistency measure aspects of the language model itself, many agentic applications depend on components like a retrieval system (as used in RAG pipelines). For these applications, evaluating the quality of the retrieval step is a critical part of assessing the overall system. This component evaluation ensures that the LLM is provided with relevant, high-quality context, which is essential for generating a useful and accurate final response. This focus

on the retrieved context makes it distinct from a metric like consistency, which measures the logical coherence of the final generated output.

A system designed to measure retrieval quality should evaluate both the relevance of retrieved documents and their effectiveness in supporting the system's overall goals, such as providing accurate answers. It should track how often the most useful information appears at the top of search results and identify gaps where relevant content is missed. To achieve this, the system should combine automated metrics with human feedback. Additionally, it should support continuous improvement by highlighting patterns in errors, enabling adjustments to search algorithms, ranking methods, and data quality to enhance retrieval accuracy over time.

By measuring and monitoring retrieval quality, developers of agentic applications can ensure that their systems consistently deliver accurate, relevant information, leading to more reliable and effective user interactions. This process helps identify weaknesses in how information is retrieved and ranked, allowing for targeted improvements in search algorithms and data management. Ultimately, it supports the development of smarter, more responsive applications that can better understand and meet user needs.

## Extrinsic Quality Evaluation

What does it mean for the output of an agentic application to be good for a particular purpose? Given the broad range of potential contexts, the best we can say is that a good output is one that serves its intended purpose well. This means that, in order to measure the quality of the output, you will necessarily need to define what good looks like on a per-application basis, often with input from human experts. For some tasks, this definition can be encoded into a “golden dataset” to serve as ground truth for automated evaluation.

However, creating such a dataset is a nontrivial process, and for many open-ended agentic tasks, it may not be feasible at all. Furthermore, since agreement among experts is rarely perfect, any golden dataset inevitably reflects subjective choices rather than a single objective truth. These are the fundamental challenges of aligning the behavior of an agentic application with real-world business objectives.

Take, for example, an application that categorizes customer service requests according to a company's specific product line. For the application to work properly, the request must return only the name of the category, not an entire sentence. In such an application, what constitutes a good answer is two-fold:

*The accuracy of the response*

Did the LLM correctly classify the request according to the organization's categories? This would likely require a golden dataset to test against where, for example, customer service requests are categorized correctly, given that the knowledge is specific to the organization and may be difficult for an LLM to determine without this additional guidance.

*The format of the response*

Did the LLM respond with only the category title, as instructed? This can be evaluated simply with a small glossary and a predefined rule.

Now imagine an agentic application where the agent must choose the appropriate tool for a task from among a list provided to it and then interact with that tool correctly based on the information in the schema. This will require a different test to ensure that the task is being completed correctly.

Once you have defined what a desirable output is for your agentic application, you then must define a metric for that measurement and a method to generate that metric. There are many methods available, and when implementing an LLM Mesh, it is important to ensure that you make a wide range of methods available to your application developers so that they can choose the right combination of methods and metrics for each application. In this way, the LLM Mesh will save developers time by ensuring that they don't have to waste time on implementing methods for quantifying the quality of their applications and ensure that the organization can be confident in the results.

There are three main categories of techniques for measuring the quality of the output of LLMs, each with its own advantages and disadvantages. Those categories are human expert methods, statistical methods that compare responses to ground truth, and LLM-based evaluation methods. The following sections review those categories.

## Human expert methods

The most reliable—but least scalable—method for measuring the quality of the output of an LLM is simply asking a knowledgeable person if the response is good or not. This approach is well-adapted during the development phase of a new agentic application. For example, in a data analytics-generating application, a data analyst could provide input on how they would solve a given problem, and the LLM could be prompted to deliver similar results.

However, once an application is deployed, the volume of content generated would make it infeasible to use human evaluation to monitor the quality—you are not going to pay a human to review every output of the LLM. That said, there are two ways that human feedback can be used to monitor quality once the application is deployed:

### *User sentiment tracking*

Include a method for simple user sentiment tracking using a binary feedback button (e.g., thumbs-up, thumbs-down). Given the low effort required of the user, it is possible to collect a meaningful sample of responses, though they may lack detailed contextual feedback on why the user responded the way they did.

### *Expert review*

Have a human expert analyze a sample of results. Such checks are an important part of quality monitoring, and they should be designed to ensure that the expert user reviews a representative sample of all responses.

When using human evaluation, it is important to ask the human evaluator to rate the output in a consistent manner across several dimensions. Simply asking if the response is “good” will not get you the information that you need to improve the application. Depending on the context of the application, you may consider asking the experts to evaluate the responses across some of the following dimensions:

### *Relevance*

Does the output align well with the query and user intent? How well does it address the core needs of the business use case?

### *Accuracy*

Are the facts and information presented correct and free of errors? This is especially crucial for applications in fields requiring high precision, like finance or healthcare.

### *Clarity and coherence*

Is the output easy to understand, logical, and well-structured? A high-quality response should be clear and devoid of ambiguous or confusing language.

### *Completeness*

Does the output provide a sufficiently complete answer, or are key details missing? Depending on the use case, an answer that is too brief or superficial may not be helpful.

### *Conciseness*

Is the response free of unnecessary information or verbose explanations? It's often important for enterprise applications to deliver only what's needed, especially when users may need quick answers.

### *Actionability*

For applications with practical implications, can the user easily take action based on the output? This is relevant in customer service, recommendation systems, or task-based LLM applications.

### *Tone and style*

Does the tone fit the enterprise's needs? For instance, customer-facing applications may need a friendly tone, while internal documentation tools might require a formal, straightforward, and technical approach.

### *Bias and fairness*

Is the response free from harmful or biased statements? Evaluating for fairness ensures inclusivity and adherence to ethical standards.

### *Safety and compliance*

Does the output avoid unsafe suggestions or violations of regulatory standards? This is essential in sensitive domains like legal, financial, or medical applications.

### *Adaptability and contextual awareness*

Can the LLM handle context changes or follow-up questions accurately? This dimension matters in dynamic environments where information evolves or multistep tasks are involved.

### *Novelty and creativity (if applicable)*

Does the LLM offer innovative solutions or ideas? This can be particularly valuable in domains like marketing or R&D, where unique insights are beneficial.

When implementing an LLM Mesh, make sure that it can capture structured evaluations from human experts. When using multiple raters, it should provide tools to measure consensus, such as multirater agreement scores. Crucially, rather than discarding disagreements, the Mesh should use disagreement sampling to identify cases where experts differ. These ambiguous cases are invaluable for refining evaluation criteria and improving the application.

## **Ground truth–based statistical methods**

For certain use cases, you will have documented examples of correct outputs that you can use as your ground truth. In these cases, you can set up quality measures that compare the model output to the ground truth. Several of these statistical methods work by comparing  $n$ -grams, which are simply contiguous sequences of  $n$  words from a text. For example, in the sentence “The quick brown fox,” the 2-grams (or bigrams) are “The quick,” “quick brown,” and “brown fox.” By measuring the overlap of these  $n$ -grams between a generated text and a reference text, these methods can quantify their similarity.

The following methods are frequently used:

### *Translation: bilingual evaluation understudy (BLEU)*

Measures the overlap of  $n$ -grams between the generated output and a reference text. It’s commonly used to evaluate the quality of machine translation.

### *Summarization: Recall-Oriented Understudy for Gisting Evaluation (ROUGE)*

Primarily used for summarization, ROUGE is focused on recall. It measures how much of the information present in a reference summary is successfully captured by the generated summary. In

essence, it answers the question, “Did the summary remember all the important points?”

### *BERTScore*

As an alternative to relying on  $n$ -grams, BERTScore uses the **BERT pretrained transformer model** to compute the similarity between the generated text and reference text on a token level. This method captures semantic similarities rather than just  $n$ -grams. It is useful in evaluating the responses in chatbot applications, as well as for evaluating the quality of translations and summarizations.

While these methods output similarity scores (like BLEU or ROUGE), these scores can be used to derive metrics familiar to machine learning developers. By setting a threshold on a score—for example, defining any response with a score above 0.8 as “correct”—you can then classify each output as a success or failure. This allows you to calculate overall accuracy, F1 score, precision, and recall for an entire dataset.

Your LLM Mesh should make these methods available to your developers, as they can provide a useful point of reference for evaluating the quality of LLM outputs. That said, most methods are appropriate only for specific cases and cannot evaluate more complex responses. They also have the benefit of requiring relatively few computational resources to compute, in contrast to the LLM-based methods described in the next section.

### **LLM-based evaluation methods**

When performing extrinsic evaluations of agentic applications, moving beyond the limitations of traditional statistical methods and human expert methods requires turning to a method that can interpret and analyze the applications’ varied and open-ended outputs: evaluation methods that use LLMs themselves. Though potentially counterintuitive, LLMs can be used to evaluate the quality of their own output when they are carefully instructed on how to do so, similar to how a human instructor grades the tests of their human students.

It should be noted that this method—often referred to as LLM-as-a-judge—is an area of active research and experimentation. The techniques described in this section are liable to evolve or be superseded by improved methods in the near future.

The core notion of LLM-as-a-judge methods is that you develop a prompt or series of prompts that will instruct an LLM to evaluate a specific response aspect. A response aspect is simply a formal definition of the different qualities that you may be looking for in a response. **Table 4-2** summarizes these aspects.

*Table 4-2. Overview of the different response aspects that you may wish to evaluate, adapted from an article on GPTScore<sup>a</sup>*

Aspect	Task	Definition
Semantic coverage	Text summarization	How many semantic content units from the reference text are covered by the generated text?
Factuality	Text summarization	Does the generated text preserve the factual statements of the source text?
Consistency	Text summarization, dialogue response generation	Is the generated text consistent in the information it provides?
Informativeness	Text summarization, data to text, dialogue response generation	How well does the generated text capture the key ideas of its source text?
Coherence	Text summarization, dialogue response generation	How much does the generated text make sense?
Relevance	Dialogue response generation, text summarization, data to text	How well is the generated text relevant to its source text?
Fluency	Dialogue response generation, text summarization, data to text, machine translation	Is the generated text well-written and grammatical?
Accuracy	Machine translation	Are there inaccuracies missing or unfactual content in the generated text?
Interest	Dialogue response generation	Is the generated text interesting?
Engagement	Dialogue response generation	Is the generated text engaging?
Specificity	Dialogue response generation	Is the generated text generic or specific to the source text?
Correctness	Dialogue response generation	Is the generated text correct or was there a misunderstanding of the source text?
Semantically appropriate	Dialogue response generation	Is the generated text semantically appropriate?

Aspect	Task	Definition
Understandability	Dialogue response generation	Is the generated text understandable?
Error recovery	Dialogue response generation	Is the system able to recover from errors?
Diversity	Dialogue response generation	Is there diversity in the system responses?
Depth	Dialogue response generation	Does the system discuss topics in depth?
Likeability	Dialogue response generation	Does the system display a likable personality?
Flexibility	Dialogue response generation	Is the system flexible and adaptable to the user and their interests?
Inquisitiveness	Dialogue response generation	Is the system inquisitive throughout the conversation?
<sup>a</sup> Jinlan Fu et al., “GPTScore: Evaluate as You Desire,” preprint, arXiv, February 13, 2023, <a href="https://oreil.ly/2ykxK">https://oreil.ly/2ykxK</a> .		

Note that in setting up an LLM-as-a-judge system, you have many choices, including:

- Which LLM should be used as the judge? Within an LLM Mesh, this can be any of the LLM services made available to the developers. Certain LLM-as-a-judge methods seek to provide good results with smaller models, reducing the cost of the evaluation.
- Can the evaluation be completed with a single interaction with the LLM (called a single-turn method), or will it require multiple interactions (called a multiturn method)?
- Does the evaluation method require reference answers or not? If the method requires reference answers, you need to develop or otherwise provide the golden dataset that the LLM judge will refer to.
- Do the evaluations of the LLM judge correlate with the responses of a human expert judge? This is the key question! When setting up an LLM-as-a-judge method, the goal is to ensure that its responses correlate with the responses of a human expert, meaning that if a human expert would rate one aspect of a response positively, then the LLM judge would rate it in a similar manner.

Note that as LLM-as-a-judge systems are based on LLMs that are, by definition, nondeterministic, the same judge may not always give the same evaluation to the same response. That said, the same can be said about human expert judges. Your goal when implementing an LLM-as-a-judge method is to design the prompts in such a way that the LLM gives consistent and reliable evaluations of the responses. Multiple libraries and templates of such evaluations have been developed and are available from both proprietary and open source providers. Examples include [OpenAI Evals](#), [Arize Phoenix](#), and [RAGAS](#).

## Implementing Evaluation Methods

When implementing an LLM Mesh in your organization, you would decide which methods to make available to your developers and provide them as a shared service. The statistical and LLM-based methods described in the previous sections are all available as open source implementations from their original authors that you could freely use in your LLM Mesh.

Alternatively, rather than creating your own implementations of these open source methods, you could choose to use a third-party service, connecting it to your LLM Mesh. Like with LLM services, an LLM Mesh architecture should allow for connecting to both self-hosted and third-party hosted evaluation services. Third-party evaluation services include those that are offered by the cloud service providers (Amazon's SageMaker Clarify, Google's Vertex Gen AI Evaluation Service, and Microsoft's Azure Machine Learning prompt flow, which includes templated evaluation flows), as well as those offered by a host of emerging startups (for example, Lyzr, Humanloop, Cognition, among others, and with more emerging every month).

So, should you build your own implementations of open source evaluation methods or purchase off-the-shelf options? Implementing these methods within your LLM Mesh can be a good choice, as it allows you to make fine-grained decisions about how the evaluation service functions, including which LLM service they use. Remember, implementing the LLM-based evaluation methods means building and deploying your own agentic applications for your developers. Just as your organization will likely choose to buy some agentic applications and build others, you face the same choice here.

From the perspective of an LLM Mesh, it is important to treat these evaluation methods in a consistent manner, regardless of whether they are provided by a third party or if you implement them yourself within your LLM Mesh. When implementing your LLM Mesh, ensure that it allows you to define the following dimensions of an evaluation:

*Response being evaluated*

Associating it with the relevant application and agent capturing the version of each; this ensures that you have traceability of your evaluations and how they evolve with different experiments

*Evaluation method and service being used*

Including any information about the version of the service being used; for example, a small change to the evaluation service may result in very different evaluation results in some cases, so taking into account the version of the evaluation service is essential

*Aspect*

For example, factuality, fluency, etc. of the response being evaluated

*Metric being calculated*

Including a precise definition and the formula for its calculation

*Value of the metric evaluated*

So that you can track the evolution of the metric over time as you continue your experimentation or monitor the application while it's in production

Note that all third-party evaluation services may not expose this information in their API, meaning that you will not be able to capture it in your LLM Mesh, making it difficult to get a complete picture of performance across your portfolio of agentic applications. You should take this into account when deciding whether you want to use a third-party service or implement your own services.

Let's make this all more concrete by looking at a simple example of how this could work for someone developing an agentic application within an LLM Mesh architecture. In this example, let's imagine that the developer is working on an agent that includes a RAG-enriched response at one point. They want to evaluate and monitor if that

response is only making claims that are backed up by the documentation that it is using for its retrieval.

First, the developer would capture the *generated response* and log its content and the metadata described previously (application, LLM service, etc.). Then, they may choose to use RAGAS as their evaluation *method*, as it is designed to work well with RAG applications. The *aspect* that they want to measure is whether the response provides appropriate context, and the specific *metric* is called context recall. **Context recall** is calculated by simply dividing the number of claims in the generated text that can be attributed to the source text by the total number of claims made in the generated text. If this value is 1, it means that all of the claims in the response are based on claims in the source, while a value of 0 would mean that none are. This *value* would be logged then as the result of the evaluation.

Then, a second LLM-as-a-judge might evaluate the quality of the written response, checking to ensure that it meets the expected tone. That evaluation would also need to capture the generated response, method, aspect, metric, and value for the tone of the response.

In this scenario, the two LLM judges measured different parts of the response at different moments in the agent's logic chain. Given the very different natures of these tests, it is necessary to use different evaluations. There could be several more evaluations required for such an agent, depending on its complexity. This multiplicity of evaluations shows the importance of providing evaluations as a shared service so that the developers can focus on creating and perfecting the logic of the agent and not be slowed down by the important but repetitive work of setting up robust evaluations.

## Implementing a Performance Architecture in an LLM Mesh

As described in the previous sections, when implementing an LLM Mesh you will make different quality assessment methods available to your app developers. They will, in turn, apply these methods at different levels within an agentic application, generating metrics that they can then monitor over time. But what are the best practices for how these different methods and metrics can be combined to ensure that the applications are performing as intended? While the state-of-the-art is a rapidly moving target, let's describe a simple performance

architecture that could be implemented as an organizational best practice for all agentic applications using an LLM Mesh.

### **LLM-level monitoring**

At the most basic level, your organization should put in place systems to consistently monitor the performance of your primary LLMs to ensure that they are delivering consistent performance. Changes in performance may occur when either the model is updated to a newer version or when the nature of the input changes, even slightly.

In the first case, model updates may result in unexpected and potentially degraded performance, especially if your team members are crafting highly specialized prompts. These prompts may depend on quirks of a particular model version, which may disappear when a model is updated. While it is true that newer versions of models generally offer improved performance, often at a lower cost, it is important to be able to monitor performance so as not to be caught unaware of changed performance in a production application.

In the second case, even when a model stays the same, it is possible that your input to the prompt may change over time without you realizing it. For example, if you are using an agentic application to process customer service requests, the content of those requests may shift as your organization introduces new products and services. This may result in changed performance of your application, requiring you to take some action to return the application to its desired level of performance.

In both of these cases, you need to design experiments where you compare the real-world results that your applications are generating with the expected, reference results. This can be done using the methods described previously, often using reference answers as your point of comparison. Note that if your monitoring shows that the nature of your input to the LLM has changed, it may mean that you need to update your human-approved reference answers.

### **Agent-level monitoring**

The output of an agent may vary depending on many different factors: changes in underlying LLMs, changes in the user inputs, changes in connected retrieval services, changes in tools, etc. Just as the power of agents is their open-endedness and the diversity of

systems that they can integrate, so too is this diversity a source of difficulty when trying to monitor agents' performance and diagnose any issue.

At the heart of agent-level monitoring is measuring if each step in an agent's execution is driving the ultimate task it is meant to accomplish.<sup>1</sup> To do this effectively, you need a clear set of evaluation criteria or, for test cases, a ground-truth dataset. Answering questions such as "Did the agent choose the right tool for this task?" or "Was the final output in the appropriate tone?" requires comparing the agent's behavior against a predefined definition of what is "correct" in that context.

Whether you should measure the quality of every task completion or just a sample depends on several factors:

- The criticality or riskiness of the application. Some applications are critical to core business processes, and others are high risk (and some are both!). In these cases, you may choose to monitor all completions, rather than a sample.
- The variability of the performance. If the agent shows that it can complete the task in a consistent way, then you may content yourself to measure the performance of only a sample.
- The volume of task completion and cost of the monitoring. As described previously, certain evaluation methods can be costly in and of themselves. Some agents may be completing many thousands of tasks per day. You will need to use the cost-measuring techniques in [Chapter 3](#) to fully understand and capture the costs of these processes to ensure that you have the budget for them.

### Agent self-monitoring

You will recall that in the previous sections of this chapter, we described LLM-as-a-judge methods of quality assessment as being essentially agentic applications themselves. You also recall in [Chapter 2](#) that we said that agents could call other agents as tools. So, agents can also use quality assessment agents to monitor their own

---

<sup>1</sup> We see again here yet another reason why it is important to design agents with a relatively narrow scope of action. When the scope is broader, it is more difficult to specify what good task completion looks like.

outputs. Agents can use quality assessment tools to monitor and attempt to correct their own outputs within a single execution run. For example, an agent might check its own draft response and, if the quality is low, try again.

This self-correction process should be fully logged. In practice, these logs provide the data for long-term iterative improvement, which still relies heavily on offline, human-in-the-loop processes. By reviewing the logs of these self-correction attempts, developers can identify common failure modes and improve the agent's design by refining prompts or updating tool schemas. As always, reference answers and human expert monitoring remain essential to ensure high quality and adherence to ethical and regulatory requirements.

## Measuring and Monitoring the Speed of Agentic Applications

Like any other service, it is important to monitor the speed and responsiveness of the agentic applications that you will build. As discussed in [Chapter 2](#), an agentic application involves many API calls to LLM services and other objects within the LLM Mesh. The speed and responsiveness of the agentic application will depend on the collective speed and responsiveness of these various services. Thankfully, monitoring the speed of API services is a well-established DevOps practice, and those methods apply equally well here. Thus, this guide will not develop those concepts fully but rather only mention a few aspects that are specific to agentic applications.

The speed of agentic applications is measured using specific metrics, primarily *end-to-end latency* and *throughput*. For an agent, end-to-end latency is the total time from initial user input to the final response. It is crucial to recognize that this includes not just the LLM's inference time but also other significant factors like network overhead, the time required for external API calls, and the execution time of any tools the agent invokes.

The LLM-specific component of latency, known as inference latency, is usually measured in time to first token (TTFT) and time per output token (TPOT). These can be used to calculate the total generation time for a single LLM call.

Inference throughput is measured primarily with tokens per second. This metric most often takes into account only output tokens, and an LLM Mesh should specify whether this is the case or not.

An LLM Mesh should provide observability of these performance metrics to ensure that they are meeting the requirements of an application. These requirements should be documented in the LLM Mesh on a per-application basis as well. For example, an internal chatbot supporting customer support agents as they work in real-time with customers will have different and higher performance requirements than an application that runs silently in the background analyzing contracts with suppliers.

## Capturing and Optimizing the Costs of Performance Evaluation

As you have probably already realized, LLM-based evaluation methods can result in a lot of traffic to your LLM services and will generate costs. It is important to capture and optimize these costs as you would for any agentic application, as described in [Chapter 3](#).

A key part of this is being intelligent about the frequency and scope of evaluations. For production applications, evaluating every response is rarely feasible due to cost. Instead of simple random sampling, more sophisticated methods can provide better insights, such as stratified sampling to ensure coverage across different interaction types, anomaly detection to focus on unusual outputs, or active sampling to select cases where the model is least confident. This decision must be balanced against the criticality of the application.

Furthermore, using large, proprietary models for LLM-as-a-judge evaluations can become expensive quickly, potentially making an application economically unviable. A more cost-effective strategy is to use smaller, open source models that have been fine-tuned for specific evaluation tasks. An LLM Mesh is the ideal architecture for this approach, as it allows you to easily integrate and manage a fleet of specialized, fine-tuned evaluator models, giving developers the flexibility to choose the right tool for each evaluation.

# Conclusion

In this and the previous chapters, we have seen how we can quantify and reduce costs, as well as how to measure performance in terms of both the quality of the response and the overall speed of the service. By measuring these three dimensions, your developers will be able to experiment with different approaches to find the right combination of speed, quality, and cost. Once again, there is no silver bullet and no single best practice. But an LLM Mesh can make it far more efficient for teams to test and build high-performing, cost-optimized applications by ensuring that they focus on the development of the application itself and not the supporting evaluation and monitoring services since those are provided by the LLM Mesh.

In the following two chapters, we will discuss additional shared services that you should also include when implementing an LLM Mesh. These are services to ensure the safety and appropriateness of the content generated by your agentic applications. As you will see, this is a subset of the quality measurement topics covered in this chapter. Then, we will address the security requirements for running agentic applications in the enterprise. As anyone who has built and attempted to deploy an enterprise application can attest, if an application cannot satisfy security requirements, there is no way it can ever be deployed. As with cost, quality, and safety monitoring, an LLM Mesh can also help here by providing the necessary mechanisms to pass those tests with flying colors.

# Safety

A “safe” agentic enterprise application is one that is reliable, harmless, and compliant. Safety and quality are closely related, but for enterprise applications, safety must be treated as a distinct and primary requirement. While quality, as discussed in [Chapter 4](#), focuses on metrics like correctness and coherence, safety ensures that an output is also harmless and compliant. For example, a response could be factually accurate yet use hateful language; it might score high on a quality metric for accuracy but would fail the safety test. In practice, this means an enterprise must treat safety as a nonnegotiable prerequisite: an answer that is insightful but violates company ethics or policies is unacceptable, regardless of its quality.

Agentic applications can expose the organization developing and deploying them to legal and reputational consequences. An unsafe output (e.g., one containing defamatory or biased content) can lead to compliance violations or liability for the company. For instance, if a customer support chatbot produced a harassing or discriminatory message, the organization could face HR and legal consequences. Therefore, safety in an LLM Mesh is directly tied to compliance (adhering to laws and regulations) and liability risk management.

It is useful to understand how an LLM Mesh architecture can support your safety efforts. Some of the LLM’s contributions to safety are direct, while others are indirect.

An LLM Mesh should include certain capabilities that are specifically designed to support the development and deployment of safe agentic applications. These capabilities include:

*Modular and reusable safety services*

Just as cost and performance measurement services can be reused and combined, safety services such as content filtering and hallucination detection can be combined in a flexible, layered way.

*Transparent and centralized logging*

As a central governance and orchestration layer, an LLM Mesh provides complete, auditable logs that can be used to identify and trace any unsafe behavior.

*Technology optionality*

In the event that a provider of AI services is compromised, an LLM Mesh allows for a rapid and low-friction migration to an alternative provider.

In addition to these direct contributions, an LLM Mesh also supports the development and deployment of safe agentic applications indirectly through some of its common capabilities. Examples of safety practices supported indirectly by an LLM Mesh include:

*Adversarial and robustness testing*

An LLM Mesh can provide components used for this testing, such as prompt libraries, automated testing agents, and result analysis tools, and support the collection and analysis of results.

*Deployment and monitoring*

Safety and robustness testing needs to continue post-deployment, which can be supported by the performance and safety testing frameworks provided by an LLM Mesh.

*Grounding agentic applications in enterprise data*

By facilitating secure connections with enterprise data, an LLM Mesh can help to ensure that an agentic application's outputs are based on and informed by real-world facts.

In this chapter, we will explore the ways that an LLM Mesh can directly support the development of safe agentic applications in the enterprise. This chapter is divided into sections focusing on the following five safety capabilities for developing and deploying safe agentic applications:

*Safety filtering*

This detects and removes unsafe content using dictionary, NLP, or LLM-based techniques.

*Hallucination detection and mitigation*

This detects when an LLM may be generating content that is not based in fact and provides remediation options.

*Transparency and explainability*

This ensures that the inputs to an LLM are transparent and that the overall logical structure of an agentic application is understandable.

*Adversarial and robustness testing*

These two types of testing subject agentic applications to simulated attacks and unexpected situations to identify potential weaknesses.

*Human feedback*

Human experts must remain the ultimate arbiters of the correctness and appropriateness of the outputs and behaviors of agentic applications. An LLM Mesh can help collect human feedback in a consistent and centralized manner.

Note that we are treating “safety” and “security” as separate topics in this guide. The topic of safety concerns the content generated by or input into agentic applications, while security concerns which entities have access to which components of an LLM Mesh. Security is the topic of [Chapter 6](#).

## Unsafe Behavior in Agentic Applications

Before detailing the safety practices that an LLM Mesh supports or provides, let's quickly review the ways in which agentic applications can behave in unsafe ways. The following three sections describe the different types of risks associated with unsafe behavior.

## Behavior with Cybersecurity Risk

Generative AI and the agentic applications that it powers can expand the attack surface of organizations that deploy it.<sup>1</sup> In the context of agentic applications, specifically, the risk of *direct and indirect prompt injection*<sup>2</sup> creates an additional risk for the enterprise that needs to be defended against systematically.

## Behavior Exposing Strong Legal Liability

There are a number of ways that agentic applications can potentially expose the organization that develops and deploys them to strong legal liability. These include:

### *Toxic content*

Content that is biased, violent, hateful, racist, sexist, or sexually explicit can run afoul of both internal policies and local laws. Furthermore, even if it is not illegal, it can present a massive reputational risk to the enterprise.

### *Forbidden content*

Some content may not be toxic but is still forbidden for an agentic application to divulge. This includes personally identifiable information and IP-protected content, which could be both in text and code.

### *Expert advice*

Some types of advice are governed by strict regulations that stipulate who can give them and under what circumstances. For example, agentic applications must refrain from providing medical, legal, or financial advice.

## Unwanted Behaviors

This category addresses behaviors that do not expose the organization to strong legal liability but are nevertheless unwanted. Examples include:

---

1 Capgemini Research Institute, “Generative AI in Cybersecurity,” Capgemini, November 2024, <https://oreil.ly/tAkIu>.

2 Indirect prompt injection is where the malicious content is not directly in the prompt but in content that the prompt directs to (e.g., in the vector store).

### *Digressions*

Given the nondeterministic nature of agentic applications, it is possible for them to go beyond their intended scope of action. If, for example, an organization has developed a sales pipeline evaluation agent that can access the organization's HR systems in order to analyze the data on a per-employee basis, the agent should not read, update, or expose information from the HR systems that is not related to sales pipeline analysis.

### *Hallucinations*

This well-known phenomenon, where LLMs generate outputs that are factually unsupported, nonsensical, or contradict known information, regardless of presence in the training set, can have strongly negative consequences in the context of enterprise agentic applications, where end users may be making important business decisions based on the output of the agent. Depending on the nature of the hallucination, it could also expose the organization to strong legal liability or a cybersecurity risk.

Various resources exist to provide a more in-depth taxonomy of risks associated with LLM-powered systems, including AEGIS2.0 by NVIDIA.<sup>3</sup>

Now, let's look at each of the five critical safety capabilities provided by an LLM Mesh.

## **Layered and Reusable Safety Filters**

Safety filters are an integral component of an LLM Mesh architecture. They are embedded systematically across every layer of the system's operational stack, including:

- The LLM level
- The knowledge base retrieval process
- Tool invocation layers
- Agent-level interactions

---

<sup>3</sup> Shaona Ghosh, "AEGIS2.0: A Diverse AI Safety Dataset and Risks Taxonomy for Alignment of LLM Guardrails," preprint, arXiv, January 15, 2025, <https://oreil.ly/zUyBU>.

By providing safety filters as shared services that can be embedded into each of these layers, an LLM Mesh allows developers to build agentic applications while mitigating the risk of unsafe, harmful, or inappropriate content ever reaching the end user.

Figure 5-1 illustrates how layered and reusable safety filters can be deployed in an enterprise agentic application. Imagine that this agent partially automates the resolution of technical support tickets. At the center is the agent itself, which is surrounded by content filters, each specialized to detect a particular type of harmful content. Each of the connected components, including the two LLMs, the two tools, and the external knowledge base, are all protected by a prompt injection filter. Finally, each of the external tools and the knowledge base is also protected by a forbidden content filter. Together, this is an example of a robust safety architecture for this type of agentic application.

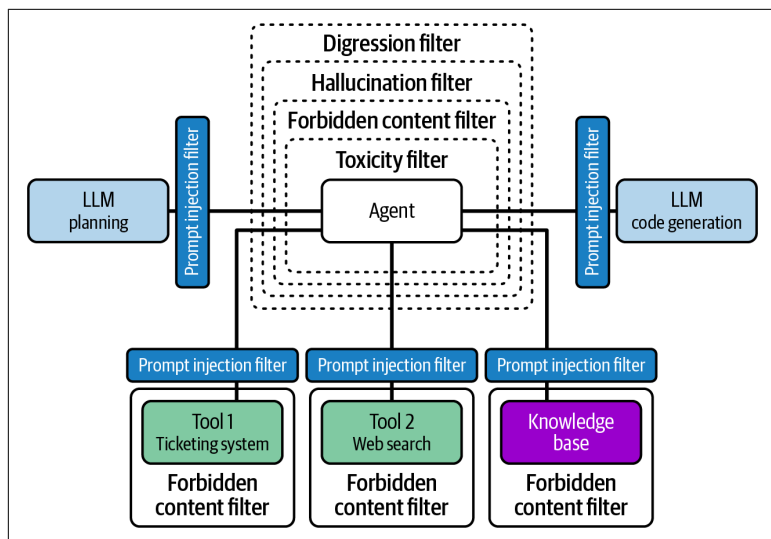


Figure 5-1. Layered safety pipeline in an example agentic application

Safety filters in the LLM Mesh are built as modular, reusable components to maximize flexibility, consistency, and operational efficiency. Each filter can be understood as a self-contained wrapper or middleware around individual LLMs or agents. This modular design ensures that safety filters can be consistently deployed across diverse contexts and multiple applications throughout an enterprise. Reusability reduces redundancy by ensuring that common

components are not needlessly redeveloped, promotes standardization of safety protocols, and significantly streamlines the effort required for developers to maintain robust, enterprise-wide safety practices.

Filters within an LLM Mesh can also operate in different behavioral modes to handle safety risks appropriately. Some filters operate in a self-correcting mode—automatically detecting and removing harmful or toxic content while allowing the underlying agent or LLM to continue its operation seamlessly. Conversely, recognizing that detection of sensitive content is often probabilistic, critical filters can operate in a fail mode. Upon detecting potential high-risk triggers, such as content flagged as possible PII or explicitly dangerous instructions, they halt the agent's response immediately and trigger a request for human review. This dual-mode flexibility ensures that the LLM Mesh can dynamically enforce both nuanced and stringent safety boundaries aligned precisely to the enterprise's specific compliance requirements and risk tolerance.

By combining filters with different behavioral modes, developers can easily create a layered safety pipeline, as illustrated in [Figure 5-1](#). Relying on a single safety filter is inherently risky, as individual filters can miss subtle threats or edge cases. A layered safety pipeline instead follows a defensive approach, strategically placing different filters at different points of the application's execution flow.

For example, within an LLM Mesh, a developer might implement initial input filtering of prompts to identify and block malicious, forbidden, or inappropriate user requests before they reach the LLM itself. Then, at the point of content generation, output filtering techniques can further inspect the generated response before it reaches the user. Each layer can employ multiple complementary methods: a keyword-based blacklist and a machine-learning-based content classifier might both be applied to incoming prompts. If either technique flags problematic content—such as forbidden terms or toxic language—the system can immediately refuse or modify the input.

Similarly, output-stage filters might combine model-based toxicity detection with simpler, rule-based checks—such as regular expressions designed to identify PII. By strategically stacking these complementary safety filters, the LLM Mesh creates multiple independent opportunities to detect unsafe content. This ensures robust protection—if one layer inadvertently allows unsafe content

through, another layer is likely to intercept it, providing resilient and comprehensive safety coverage.

## Methods for Detecting and Filtering Unsafe Content

When implementing an LLM Mesh, consider providing both third-party hosted content moderation services as well as self-hosted services based on open source frameworks to your application developers. This way, the application developers have the option to use the most appropriate service or combination of services for the given application they are using. For example, in many cases, a good solution can be to use a self-hosted open source framework as a first layer of moderation and fall back to a third-party service for a second opinion on complex or ambiguous cases.

Most content moderation services and frameworks rely on a combination of techniques to detect and classify text as safe or unsafe. These may include rules-based techniques (e.g., a dictionary of forbidden terms) and pattern-matching techniques (e.g., regular expressions to detect content formatted in a distinctive way, such as email addresses or Social Security numbers).

Examples of externally hosted content moderation and filtering services include:

- **OpenAI Moderation API**
- **Azure AI Content Safety**
- **Perspective API** (developed by Jigsaw and Google)

A variety of open access content moderation models are available as well, some of which are specific to certain types of content. These could be deployed as options available to developers in an LLM Mesh. Examples include:

- **Presidio**, a data protection and de-identification SDK from Microsoft
- **NeMo Guardrails** family of models developed by NVIDIA
- **ShieldGemma** family of models developed by Google
- **spaCy**, an open source NLP library in Python

Choosing between a self-hosted or a third-party service requires taking into account the same trade-offs for LLM services, as discussed in **Chapter 1**. While third-party hosted services provide

the convenience of zero maintenance, they come with a cost per use and the requirement of sending data outside of your organization's firewall. While self-hosted services provide fixed costs and the guarantee that data does not leave your organization, they require effort to set up and maintain.

## Hallucination Detection

Hallucinations are instances where the LLM produces information that sounds plausible but is factually incorrect or fabricated. While not “unsafe” in the offensive content sense, hallucinated outputs can mislead users or result in poor decisions. Recent research has focused on detecting these falsehoods. One notable finding by Azaria and Mitchell (2023) is that an LLM's internal state often reveals when it is fabricating information.<sup>4</sup> Specifically, when an LLM generates a statement that is not well-supported by its training data, the probability distribution of the words it chooses may show higher uncertainty or inconsistency. These subtle probabilistic cues can be measured to identify outputs that are likely untruthful, effectively creating a “lie detector” for the model.

Building on this idea, an article introduced semantic entropy probes (SEPs), lightweight models that attach to an LLM's hidden layers to estimate the “uncertainty” of a response in semantic space.<sup>5</sup> High semantic entropy correlates with hallucination; essentially, if the model's internal embeddings show it is unsure or deviating from known facts, a SEP can flag that answer as likely fabricated.

When implementing an LLM Mesh, these emerging techniques could be implemented as shared services for warning users that a response generated by an LLM may be a hallucination. This could be particularly useful for output generated in the chain of thought reasoning of an agent that is otherwise invisible to the user. Automatically flagging these probable hallucinations could greatly assist the developer of an agentic application when debugging why an agent is performing improperly.

---

4 Amos Azaria and Tom Mitchell, “The Internal State of an LLM Knows When It's Lying,” preprint, arXiv, October 17, 2023, <https://oreil.ly/3tQtL>.

5 Jannik Kossen et al., “Semantic Entropy Probes: Robust and Cheap Hallucination Detection in LLMs,” preprint, arXiv, June 22, 2024, <https://oreil.ly/Pcr75>.

## RAG as a Hallucination Mitigation Strategy

Beyond detecting hallucinations after they occur, there are measures you can take to prevent a hallucination from happening in the first place. One preventative measure is RAG, which grounds an LLM's answers in factual references from an external knowledge base to reduce hallucination frequency. However, RAG is not a complete solution; its effectiveness depends on the quality of the retrieval. If the system retrieves irrelevant or incorrect documents, it can propagate errors, and the LLM may still misinterpret or inaccurately summarize even correctly retrieved information. Despite these limitations, RAG provides a critical advantage: it gives humans a straightforward way to fact-check an answer by consulting the specific source documents.

This is because the LLM doesn't have to rely purely on the statistical correlations of its training data (which might be outdated or incomplete). Instead, it can quote or summarize the retrieved facts. For example, instead of guessing a software API's behavior, a RAG system will fetch the official documentation and let the model base its answer on that. This approach is particularly relevant when making tools and their documentation available to LLMs in agents.

In an LLM Mesh, RAG can be seen as a safety-enhancing architectural choice: by designing agents to always consult a trusted data source (company manuals, policies, technical documentation, etc.), the system can limit off-the-cuff fabrications.

## Hallucination Detection Pipeline

In certain critical applications, it may be necessary to implement a hallucination detection evaluation to ensure that answers are factual and grounded in established sources. Such an evaluation could be triggered by a SEP indicating an elevated probability of a hallucination. An example multistep pipeline for hallucination safety is:

1. After the LLM generates an answer, the LLM Mesh extracts key factual claims (names, dates, stats) from the text.
2. These claims are sent to a verification module—e.g., a tool that queries a knowledge base or even a web search API—to check if they can be corroborated.

3. In parallel, the system can run an explicit hallucination classifier: for instance, using a semantic entropy probe that outputs a score for how likely the content is correct.
4. If the answer is found unsupported or the truthfulness score is low (beyond a threshold that would be defined for the specific use case), the LLM Mesh would identify it as a probable hallucination. It could then either warn the user (“⚠ I’m not confident in this answer’s accuracy”) or trigger a fallback, such as switching to a retrieval-based answer or providing a confidence score. For example, an agentic application in a healthcare setting that answers a medical query might generate an unsupported treatment suggestion—the verification step fails to find that suggestion in trusted medical sources, so the system replaces or annotates the answer. This kind of pipeline ensures that factual correctness is assessed and enforced as part of safety, rather than assuming every LLM output is accurate.

## Explainability and Transparency

LLMs are inherently opaque, making it fundamentally impossible to fully understand or reliably explain precisely how they arrive at a particular output. This intrinsic inscrutability arises from the complexity and scale of the neural networks powering these models, where billions of parameters interact in ways too intricate to interpret at a granular level. While research labs use complex interpretability techniques, these are not feasible for enterprises at scale. For enterprises, the pragmatic best practice is to accept this inscrutability.

In this context, explainability shifts from model introspection to system-level auditability. This includes clear logging of prompt history, API calls, decision chains, and filtering logic. The goal is for an enterprise to be able to reconstruct why an agent took a certain action, even if they cannot explain exactly why the model generated a particular token. An LLM Mesh supports both explainability and transparency, as summarized in [Table 5-1](#) and further detailed in the next sections.

Table 5-1. Interpretability, explainability, and transparency in an LLM Mesh

	Definition	How it is supported by an LLM Mesh
<b>Interpretability</b>	The ability to understand the internal mechanics of how an LLM produces a specific output	Not directly supported; the LLM Mesh instead provides robust tools for transparency and explainability
<b>Explainability</b>	The ability to explain how an LLM or multiple LLMs interact with other systems in an agentic application	Standardized, modular objects Centralized, auditable logging
<b>Transparency</b>	Fully capturing all inputs into an LLM, associated with the resulting response	Centralized, auditable logging

Enterprises must acknowledge the inherent opacity of LLM inference processes. Rather than striving to fully interpret every internal decision of the model—something currently impossible at scale—organizations should focus explainability efforts on transparent logging, auditing, and accountability around data flows, model choices, and user interactions.

## Transparency in Data Flow and Model Selection

Transparency within an LLM Mesh primarily revolves around making clear what data informed an LLM’s response and which objects of the LLM Mesh participated in its creation. This involves capturing precise details about the prompt provenance—such as whether the prompt originated directly from a human user or was generated programmatically by another agent. Similarly, when using RAG, the specific documents retrieved to provide context must be transparently logged.

Transparency also extends to documenting the choice of LLM or tool selected to fulfill the request. For instance, when responding to a query requiring complex reasoning, an application builder working in an LLM Mesh architecture might choose a large and powerful model such as GPT-4o, clearly recording this selection in its logs. Conversely, for simpler factual queries, a smaller, more cost-effective model might be selected, again with clear traceability. By explicitly documenting the chain of events, decisions, and dependencies, the LLM Mesh provides enterprises with genuine transparency without trying to expose the untraceable inner workings of the LLM itself.

## Achieving Explainability Through the LLM Mesh Architecture

Given the fundamental opacity of LLMs themselves, explainability within the LLM Mesh is instead achieved by making transparent every step that surrounds the call to the LLM service. Specifically, the LLM Mesh's architecture allows enterprises to track and explain the flow of data from the initial prompt to the final delivery of the response.

By explicitly logging each external interaction—such as user inputs, agent choices, retrieval decisions, and filtering actions—an LLM Mesh provides comprehensive traceability. While logging does not in itself guarantee explainability, this traceability is the foundation for auditability. The quality of this auditability, however, depends entirely on the granularity and consistency of the logs. This systematic documentation allows stakeholders to reconstruct the sequence of events and understand which data and components influenced the generation process, providing a practical alternative to interpreting the opaque internal state of the LLM.

## How an LLM Mesh Supports Explainability and Transparency

Without the centralized control and explicit structure provided by an LLM Mesh, ensuring explainability and transparency can quickly become a daunting task. Imagine an organization that is developing agentic applications without an LLM Mesh architecture. Multiple teams would use different LLM services independently, resulting in inconsistent documentation, logging formats, and transparency standards. Consequently, tracing the data provenance, understanding the reasoning behind model selection, or auditing specific interactions would require painstaking manual effort, potentially becoming impossible.

By contrast, an LLM Mesh architecture inherently standardizes and centralizes these processes, systematically capturing details like prompt provenance, retrieval decisions, model selection, and filtering outcomes. This centralized transparency significantly simplifies regulatory compliance and ethical oversight, ensuring that every AI interaction is clearly documented, consistent, and easily auditable—even though the internal decision making of LLMs themselves remains inherently inscrutable.

# Adversarial and Robustness Testing

Just as software undergoes penetration testing, agentic applications that are used in exposed contexts, such as a chatbot open to the public on a corporate website, or sensitive contexts require adversarial testing, frequently called *red teaming*. This technique uses carefully crafted prompts and scenarios to deliberately bypass an agentic application’s safety guardrails.<sup>6</sup> Robustness testing seeks to test the application’s ability to continue to operate safely even in rare or unexpected conditions. The following two sections describe these in more detail, and a final section on continuous monitoring explains that these practices must be continued post-deployment.

## Adversarial Testing

Adversarial testing involves creating input prompts designed to jail-break the model or expose biases. Techniques range from simple role-play prompts (“Pretend you are an evil AI...”) and encoding attacks that hide malicious instructions from filters to exploiting the agent’s function-calling capabilities by tricking it into misusing a tool. While this testing is critical for identifying vulnerabilities, it is also a resource-intensive and ongoing process. No suite of tests can be fully comprehensive given the evolving nature of attacks. Therefore, the goal is continuous mitigation rather than perfect prevention. Frameworks like HarmBench can help standardize this process by providing structured scenarios for evaluation.

## Robustness Testing

Robustness testing seeks to ensure safety under stressful or novel conditions. This includes testing the agentic application with unexpected inputs—such as irrelevant or malicious code, random Unicode characters, or contradictory instructions—and monitoring for undesirable changes in behavior. For applications relying on proprietary models, it also means regularly monitoring for model degradation, as unannounced provider updates can unexpectedly alter performance and behavior. Furthermore, it involves evaluating how well the underlying models generalize to tasks, modalities, and

---

<sup>6</sup> Microsoft offers a [detailed guide](#) for red teaming LLMs and applications built using them.

languages they have not explicitly been trained on, a key challenge highlighted by recent research.<sup>7</sup> The goal is to ensure the application remains reliable and safe even when encountering situations outside of its expected operating parameters.

For example, an excessively long input could cause a content filter to time out and skip, which could be exploited. The goal of such testing is that, no matter how bizarre or tricky the input, the worst the system will do is err on the side of caution and refuse to respond, rather than produce something dangerous. Robustness testing complements adversarial testing: together they ensure the system is both hardened against known attack patterns and generally resilient to unexpected inputs.

## Continuous Monitoring for Attacks

Enterprises should integrate adversarial and robustness testing not just as a one-time activity before deployment but as a continuous monitoring and improvement loop. For instance, the teams developing agentic applications can maintain a library of “red team” prompts and include them in an agentic application designed for semi-autonomous testing of other agentic applications within an LLM Mesh. Any noncompliant answers or safety bypasses are flagged as failures to be fixed (much like unit tests).

Additionally, real-world usage should be monitored for signs of attack: if suddenly a user query pattern looks like an attempted jailbreak or a flood of queries is trying to find a policy gap, the system should alert administrators.

Adversarial and robustness testing should be thought of as a continual process where the learnings from the testing and observations are analyzed and used to improve the test procedures. **Figure 5-2** illustrates this process.

---

7 Lingjiao Chen et al., “How Is ChatGPT’s Behavior Changing Over Time?,” preprint, arXiv, October 31, 2023, <https://oreil.ly/QE90S>.

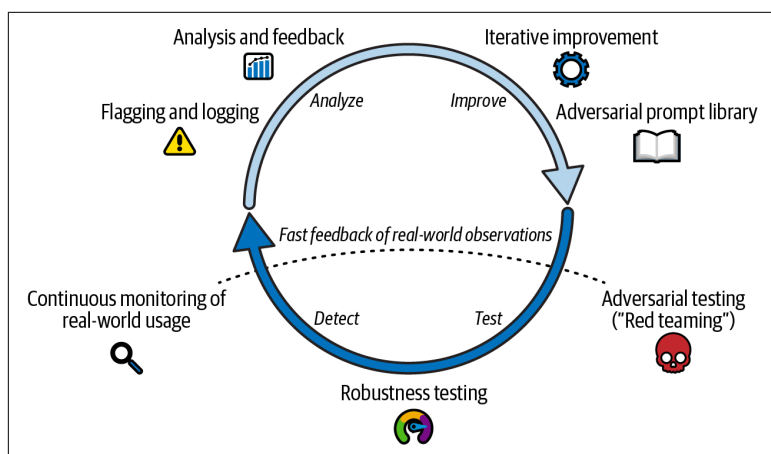


Figure 5-2. Adversarial and robustness testing and improvement cycle for agentic applications

## Human Feedback as the Ultimate Arbiter of Safety

As discussed in [Chapter 4](#), no automated system can catch every nuance of safety. Thus, human feedback and oversight are the final backstop for ensuring content is safe. In practice, human-in-the-loop (HITL) for agentic applications means **certain multistep processes are interrupted** and get routed to a human moderator or reviewer for approval.

This is especially important for edge cases and high-risk scenarios. This workflow is a direct implementation of the filter “fail mode” discussed earlier in this chapter. For example, if a safety evaluation returns an elevated level of uncertainty about the riskiness of a piece of content (say the language is borderline harassing or it’s medical advice that could be high-stakes), it can flag the response as requiring human review. A content moderator or domain expert then examines it and either approves it, edits it, or blocks it.

Human feedback is also used to improve the system over time, though the methods vary in complexity. Most enterprise systems implement simple human-in-the-loop mechanisms like escalation queues or feedback forms (e.g., thumbs-up/down ratings). This data is invaluable for identifying failure patterns and manually improving an agent’s prompts. A more complex application of this concept

is reinforcement learning from human feedback (RLHF), which foundation model developers use to align models like ChatGPT with human preferences. However, implementing a full, in-house RLHF training loop is a nontrivial undertaking and is rarely done in enterprise applications.

Let's look at an example workflow for human moderation/review. Imagine an agentic application deployed in a customer support setting. The workflow could be:

1. The user's query comes in.
2. The LLM drafts an answer.
3. If the answer contains anything uncertain (like it had to use an internal policy document that it's not 100% confident about, or the safety classifier gave a mild toxicity warning), then instead of directly sending it to the customer, the draft is forwarded to a human support agent.
4. The human customer support technician sees both the user question and the AI's draft answer, and in a special interface, they can accept it, edit it, or provide their own answer.
5. Only after this human check does the answer go out to the customer.

Over time, the system can learn from the human's decisions (e.g., if the human technician always fixes the AI's tone in certain types of responses, that feedback can be used to improve the LLM's tone in the future).

In less critical use cases, the human-in-the-loop might be invoked asynchronously or selectively—for instance, an HR chatbot answering employees might normally answer automatically, but if an employee asks something potentially sensitive (“I’m experiencing harassment; what should I do?”), the system might escalate that conversation to an HR representative rather than have the chatbot answer. The key point is that human judgment is the fail-safe: it can handle complexities of context, ethics, and liability that an AI may not fully grasp.

## Conclusion

While the inherent complexity and opacity of LLMs present unique challenges, enterprises can nonetheless implement practical measures to manage risks effectively. Ensuring the safety of agentic applications powered by LLMs demands a multifaceted approach, combining proactive content filtering, robust hallucination detection, structured transparency, and deliberate human oversight. By embracing a layered pipeline approach, maintaining transparency at every practical interaction point, and aligning closely with evolving regulatory and ethical standards, organizations can confidently deploy LLMs without compromising their values, reputation, or compliance obligations. The next chapter explores how to complement these safety practices with rigorous security strategies, ensuring comprehensive protection of your LLM Mesh.

# Agentic Application Security

Agentic applications built with LLMs present fundamentally different security challenges compared to traditional, monolithic software. These agentic applications often consist of numerous tools, services, and models—becoming dynamic, multiobject systems that can be difficult to govern without architectural support.

Security policies are only effective if they can be applied consistently across all services and objects. In traditional architectures, this consistency is hard to achieve—each service or application may interpret and implement policy logic differently, leading to enforcement gaps and governance drift. An LLM Mesh addresses this by enabling centralized policy definition coupled with distributed, system-wide enforcement.

This chapter explores how adopting an LLM Mesh architecture provides the necessary enforcement infrastructure to secure these complex systems. It introduces the principles of system-wide access control, federated identity, secure communication gateways, dynamic permissions, audit logging, anomaly detection, and deployment isolation—all of which are required to establish a secure foundation for agentic applications. When these principles are embedded and enforced consistently across components, enterprises gain the confidence to develop, operate, and scale agentic systems securely.

## A Paradigm Shift: Security in an LLM Mesh

Agentic applications built using a monolithic architecture are often secured using bespoke, application-specific logic: access controls are tightly coupled to internal logic, services may lack standardized logging, and secrets such as API keys or credentials are often embedded directly into code. These approaches result in brittle security postures that are difficult to audit, maintain, or adapt as requirements evolve.

By contrast, an LLM Mesh architecture embraces a modular design in which security is consistently enforced through common modules upon which all applications can be built. Mesh-native components—such as agents, tools, retrievers, and LLMs—can be reused across applications, and their interactions are governed by shared policy frameworks, centralized auditing mechanisms, and consistent access controls. This enables visibility and control at a level that is infeasible in monolithic architectures.

Importantly, the need for this shift stems from the inherently integrated nature of agentic applications. Unlike standalone apps, agentic systems frequently invoke external services, interact with shared data sources, and compose multiple agents into workflows. These multiobject, multiboundary patterns expand the security surface and demand architectural controls that go beyond perimeter defenses, such as firewalls and network gateways.

**Table 6-1** compares the security characteristics of monolithic agentic applications versus those built within an LLM Mesh architecture.

The following sections will explore each security attribute in more detail. The chapter ends with a short section on aligning agentic applications with enterprise security standards.

Table 6-1. Comparison of security attributes in monolithic and LLM Mesh architectures for agentic applications

Security attribute	Monolithic architecture	LLM Mesh architecture
Access control	Hardcoded, app-specific logic with minimal reuse	Fine-grained, attribute-based controls enforced uniformly across all Mesh-native objects
Secure gateway and API architecture	Perimeter-based APIs with minimal request context validation	LLM-aware gateway that inspects prompts, validates metadata, and enforces routing policies
Permissions and secrets management	Permissions scoped to users; secrets often hardcoded	Object-scoped attribute-based access control (ABAC) policies and centralized secrets managers with rotation and isolation
Audit logging and traceability	Fragmented or inconsistent logging; difficult to correlate across systems	Structured, end-to-end logs with request IDs, object metadata, and tamper resistance
Anomaly detection and monitoring	Static rules and thresholds; limited behavioral insight	Real-time, context-aware detection of agent misuse integrated with security information and event management (SIEM) tooling
Secure deployment methods	Varies by app; isolation and validation often ad hoc	Consistent enforcement of hosting policies, resource isolation, and output validation

## Access Control Frameworks

Access control is a foundational element of enterprise security—but in the context of agentic applications built within an LLM Mesh, it must go far beyond traditional user-based authorization. This section introduces three frameworks for access control in an LLM Mesh:

- Fine-grained identity enforcement across all objects
- The shift from RBAC to ABAC models
- The need for federated identity in distributed, multiagent environments

## Fine-Grained Identity and Access Enforcement

In contrast to traditional user-centric security models used in monolithic or standalone applications, identity and access control in an LLM Mesh architecture must be reimagined. Fine-grained enforcement means that permissions are not just defined at the level of users, but extended to encompass all of the objects in an LLM Mesh, including agents, services, and tools. Each of these objects operates differently and requires context-specific control.

To illustrate, consider an enterprise in which both its analytics and HR departments are deploying agentic applications. A fine-grained policy would ensure that an agent developed by an analytics team, even when operated by an authorized user, cannot access prompts, tools, or datasets associated with the HR domain. This separation is enforced at the Mesh layer, where every interaction is mediated by metadata-rich policy checks—taking into account not only the actor’s identity but also the object type, purpose, and level of sensitivity of the application.

Such precision enables a *least privilege* model at scale, supporting both compliance and risk management goals. Least privilege refers to a security principle in which every user, agent, service, or object in an LLM Mesh is granted only the minimum access necessary to perform its function—and nothing more.

## RBAC Versus ABAC in an LLM Mesh

Access control in enterprise systems has traditionally relied on *role-based access control* (RBAC), in which users are assigned specific roles, and roles are associated with permissions. This approach is relatively easy to manage at scale, especially when roles map cleanly to organizational hierarchies or job functions. However, it lacks the flexibility to respond to dynamic conditions—such as the sensitivity of the data involved, the context of access, or the identity of nonhuman actors like agents or tools.

*Attribute-based access control* (ABAC), by contrast, determines access based on a rich set of attributes associated with the subject (e.g., user, agent), object (e.g., tool, dataset), and environment (e.g., time, location, network). This model supports more granular and context-aware decisions. For example, an agent might be allowed

to access a tool only if the request occurs during business hours, originates from a trusted subnet, and involves nonsensitive data.

In an LLM Mesh architecture—where autonomous agents interact with numerous services, datasets, and tools—ABAC offers several critical advantages. It enables:

- *Context-sensitive permissions* that adapt to real-world conditions
- *Fine-grained control* over every interaction between LLM Mesh objects
- *Reusability* of access policies across components regardless of implementation details

While ABAC is particularly effective at enforcing least privilege in complex, distributed environments, its implementation requires strong data governance. Organizations must manage the challenges of “attribute explosion,” ensure the integrity and freshness of attributes from multiple sources, and resolve any inconsistencies.

## Federated Identity Across an LLM Mesh

As agentic applications become more distributed and modular, identity federation becomes a foundational requirement for secure operations. In traditional systems, identity is often managed within the bounds of a single application or service. However, in an LLM Mesh, agents, tools, and services are frequently composed across organizational units and infrastructure boundaries.

Federated identity ensures that each actor in an LLM Mesh—whether human or nonhuman—can be consistently authenticated and authorized across all components. Federated identity integrates with enterprise-wide single sign-on (SSO) systems and Identity and Access Management (IAM) platforms, allowing credentials to propagate seamlessly as agents invoke other agents, tools, or APIs. This propagation provides the trusted identity attributes required to enforce the fine-grained ABAC policies discussed previously, making the two concepts most secure when used together.

Maintaining identity continuity and traceability is essential for enforcing fine-grained access controls, generating meaningful audit logs, and enabling security policies to follow users and services wherever they operate. Without federation, these capabilities break

down, leading to inconsistent enforcement, orphaned permissions, and potential security blind spots.

Therefore, federated identity is not a localized configuration concern; it must be enforced uniformly across the entire LLM Mesh. This ensures that all policy enforcement, auditing, and monitoring mechanisms are operating with a shared understanding of who or what is acting and under what authority.

## Secure Gateway and API Architecture

As seen in the previous section, maintaining consistent identity and access enforcement is essential across an LLM Mesh. But identity alone is not sufficient—security must also be enforced at every point where components communicate. This is where gateways and APIs come into focus as key enforcement boundaries.

### Securing Communication Boundaries in an LLM Mesh

One of the most critical aspects of securing an LLM Mesh is enforcing consistent policies at every communication boundary. This goes beyond simply protecting the interfaces to language models—it includes safeguarding the entire network of objects in an agentic application.

Each interaction across these components presents a potential attack vector: prompts passed between agents, tool invocations that trigger external APIs, or payloads delivered to an LLM service. Without well-defined enforcement points, it becomes impossible to control or monitor what flows across the system.

One of the core benefits of an LLM Mesh is that it provides a unified gateway to all objects needed for building agentic applications—including agents, tools, LLMs, and retrieval services—in a way that is both consistent and tailored to their many integrations. Rather than relying on ad hoc enforcement mechanisms between components, an LLM Mesh standardizes policy enforcement through a purpose-built gateway layer.

This gateway layer serves as a set of security enforcement points by inspecting traffic, validating metadata headers, authenticating requests, and applying routing rules based on rich contextual information. Unlike a generic API gateway, an LLM Mesh gateway is capable of parsing prompt structures, enforcing input/output

schemas, and recognizing metadata about the actor, object type, and workflow context. For instance, a gateway's logic would first authenticate the requesting agent via its token. It would then fetch the agent's security policies to verify that the request is being routed to an allowed model endpoint. Finally, it would inspect the prompt itself to ensure it does not contain sensitive data before forwarding the request. If any of these checks fail, the gateway would reject the request.

In this way, an LLM Mesh becomes the operational foundation for enforcing security policies uniformly, ensuring that all interactions between distributed objects are secure, validated, and traceable.

## Secure API Key Management

API keys and secrets are a common point of failure in application security—and in an LLM Mesh, where multiple agents, tools, and services are interconnected, the consequences of poor key management are magnified. Hardcoding secrets into source code or passing them unsafely between services can lead to unauthorized access, data leakage, or compromise of critical infrastructure. For example, a developer might accidentally commit a configuration file with an API key to a public code repository, instantly exposing it to automated scanners and attackers.

To mitigate this risk, when implementing an LLM Mesh, incorporate best practices such as the use of secure secrets managers (e.g., AWS Secrets Manager, Azure Key Vault, GCP Secret Manager), enforcement of access policies around secret usage, and automated key rotation. These practices are not unique to LLM systems, but their consistent application across all objects—including prompt routers, retrievers, orchestration tools, and agent endpoints—is especially important given the many external connections typical of agentic applications.

What distinguishes an LLM Mesh is its ability to embed and enforce these practices systematically across distributed objects. Rather than relying on each component to manage secrets independently, an LLM Mesh can offer centralized enforcement points and metadata-aware security layers that ensure secrets are handled consistently and securely throughout the architecture.

# Permissions Management

Permissions management in an LLM Mesh requires more than just assigning access rights to users. This section explores two key dimensions: first, how access is controlled between different types of objects within the LLM Mesh; and second, how these permissions can be expressed and evaluated dynamically based on context.

## Controlling Data Access at the Object Level

One of the most important aspects of securing an LLM Mesh is controlling how different objects access each other—not just users accessing services, but agents accessing data and tools, and interacting with users. Because agentic applications are composed of many interconnected services and tools, enforcing object-level access is essential to avoid overreach, reduce attack surface, and preserve data boundaries.

For example, an organization might block a tool from accessing a particular set of embeddings that contain sensitive information. This is enforced at the LLM Mesh gateway, which checks for a policy denying the query action to a specific tool-`id` if the target vector-store has a `sensitivity-level` attribute of `High`. By enforcing such attribute-based permissions at the object level, policies can be more precisely scoped and aligned with security requirements.

## Expressing and Enforcing Dynamic Permissions

Permissions in an LLM Mesh must be dynamic and context-sensitive to reflect the fluid, multiobject workflows of agentic applications. While static permissions may suffice in simple systems, they quickly become brittle and error-prone in environments where agents, tools, and services operate across organizational and temporal boundaries.

To address this, an LLM Mesh should support policies defined in machine-readable formats—such as JSON or policy-as-code systems like Open Policy Agent (OPA)—and evaluate them at runtime using object-level metadata. These policies can include conditions based on time, object classification, purpose, or the identity of the interacting agent.

This dynamic evaluation aligns directly with the principles of ABAC, where permissions are granted not based on fixed roles, but

on the attributes of the subject, object, and environment. In this way, an LLM Mesh doesn't just enable dynamic permissions—it requires them—to ensure that only appropriately scoped access is allowed at every interaction point. By operationalizing ABAC within an LLM Mesh, enterprises can enforce a consistent, least-privilege agentic architecture for the entire organization.

## Audit Logging and Traceability

Robust audit logging is essential for any secure architecture, but especially so in an LLM Mesh where workflows span multiple interacting objects. This section explores how audit logging in an LLM Mesh enables end-to-end visibility, enforces structured metadata standards, and supports forensic readiness and regulatory compliance.

### Structured Logging and Metadata Standards

Structured logs are the foundation of visibility in distributed systems—and in an LLM Mesh, they become the enforcement mechanism that enables reliable security and compliance. As introduced in [Chapter 5](#), auditable logging helps track agent behavior for safety and transparency. In this chapter, we extend that principle: structured audit logs are required to build a provable, end-to-end account of all security-relevant activity across the architecture.

To be considered structured, each audit log entry must include consistent identifiers—such as request IDs, session/user IDs, and object versions—and be captured in machine-readable formats like JSON Lines (JSONL). These formats support structured metadata, enabling filtering, analysis, and correlation across distributed services. Audit logs must reflect not only what occurred but also the surrounding context: who initiated the action, which objects were involved, and what policy boundaries were evaluated.

Without such a framework, addressing challenges like inconsistent schemas, high log volumes, and timestamp synchronization in distributed environments creates significant operational overhead. A well-implemented LLM Mesh can mitigate these issues by enforcing a unified logging schema and consistent metadata propagation across all components. By centralizing the enforcement of these best practices, the Mesh reduces operational complexity and ensures the audit trail remains reliable, manageable, and secure by design.

## End-to-End Visibility Across Agent Workflows

Configuring an LLM Mesh such that it provides comprehensive audit logging means capturing every interaction across the lifecycle of an agentic application for security and compliance purposes. This includes request and response pairs, tool invocations, model completions, and the use of context objects such as prompt templates and retrievers. Multiagent workflows—where one agent calls another or composes chains of tools—require continuity of context and causal lineage to maintain a complete audit trail.

Without a unifying architecture, audit logging becomes fragmented and difficult to correlate across systems. An LLM Mesh addresses this by ensuring that every interaction flows through consistent enforcement points, where structured audit logging policies are automatically applied. This architectural consistency transforms what would otherwise be disparate, incompatible log formats into a unified, searchable record of all system activity.

However, a “log everything” approach is often impractical due to high data volumes and significant compliance risks. A single request to a multistep agent can generate hundreds of megabytes of data, and storing such detailed activity may violate data residency laws or the General Data Protection Regulation (GDPR) by capturing sensitive information. An LLM Mesh mitigates these issues by enabling centralized and dynamic logging policies. For example, an administrator can configure rules for log sampling (capturing full traces for a subset of requests), log summarization (using another agent to create a concise summary), and log redaction (automatically masking PII before logs are written). This approach allows the system to balance visibility against cost and regulatory requirements.

## Forensic Readiness and Regulatory Compliance

Audit logs serve not just operational goals but legal and regulatory ones. Breach investigations require verifiable, tamper-evident records. To achieve this, audit logs must be cryptographically signed, which allows for tampering detection by verifying the integrity of the record. Logs should also be retained in immutable stores and be searchable for rapid forensic reconstruction. Structured audit logs also satisfy requirements for frameworks like ISO/IEC 27001 and GDPR—including data subject access requests (DSARs) and evidence of data handling boundaries.

By embedding these capabilities into the LLM Mesh, organizations ensure that auditability is not bolted on after the fact but built in from the beginning. Security, compliance, and forensic readiness all depend on this architectural consistency.

## Anomaly Detection and Monitoring

As agentic applications gain autonomy and complexity, traditional rule-based monitoring becomes insufficient. This section explores how an LLM Mesh enables continuous, context-aware detection of abnormal behaviors and integrates those insights into enterprise incident response processes.

### Detecting Suspicious Agent Behavior

Unlike traditional applications, agentic systems often involve decision making and dynamic orchestration across multiple objects. This creates a wider and more complex range of potential behaviors—making it essential to detect anomalies that may signal misuse, compromise, or unsafe execution patterns, while allowing unexpected but nonetheless approved behavior.

Anomalies in the use of agentic applications can take many forms. Repeated or looping prompt patterns may indicate an agent caught in a recursive loop or attempting prompt injection. Privilege escalation attempts—such as an agent repeatedly attempting to use tools beyond its authorization—are particularly dangerous in systems that rely on fine-grained ABAC enforcement because ABAC is designed to express precise, context-dependent policies. Additionally, agent chaining abuse can emerge when agents delegate tasks across services in ways that circumvent original policy boundaries. For instance, consider a scenario with two agents:

1. A finance agent that is blocked by policy from accessing the sensitive HR database
2. An HR agent that is permitted to access employee salary data in that database

A user could instruct the finance agent to “Ask the HR agent for a list of all employee salaries and return it to me.” The finance agent, in delegating the task, is not violating its own policy. The HR agent, which is allowed to access the data, fulfills the request and returns

the sensitive information. The data is then passed back to the user through the finance agent, effectively bypassing the security control that was meant to block it.

As [Chapter 5](#) introduced, adversarial testing is one way to uncover these vulnerabilities in advance. But because agentic workflows evolve at runtime, real-time monitoring must complement testing to catch behaviors as they emerge.

## Real-Time Alerts and Response

To be effective, anomaly detection in an LLM Mesh must go beyond static thresholds. It should incorporate contextual awareness—for instance, identifying request spikes that are inconsistent with normal user behavior or flagging an agent operating outside its expected domain.

These patterns should trigger alerts that can be routed to centralized security teams and Security Information and Event Management (SIEM) platforms. However, a key challenge is managing false positives to prevent alert fatigue. Legitimate but infrequent activities, such as an agent idling or security tests like red teaming, can be incorrectly flagged. An LLM Mesh should therefore allow administrators to configure policy exceptions, such as placing an agent in a temporary “testing mode” to suppress alerts during an authorized security exercise. This ability to contextualize exceptions is crucial for building a reliable detection system.

Because agentic applications span many services and domains, only a unified architecture—like an LLM Mesh—can consistently apply behavioral policies, instrument all object interactions, and route anomalies to enforcement and response systems in real time.

## Secure Deployment Methods

This section examines the implications of how and where LLMs are hosted, as well as how different deployment patterns affect isolation, control, and trust boundaries. It explores the security implications of cloud, on-premises, and hybrid deployments; the importance of isolating objects in shared environments; and the risks of multiagent interaction without proper output validation.

## Secure LLM Hosting Models

**Chapter 1** introduced the different hosting models available for LLMs—including cloud-managed APIs, on-premises deployments, and hybrid configurations. From a security perspective, these choices significantly influence the boundaries and trust assumptions within an LLM Mesh.

In cloud-hosted deployments, the primary concern is securing data in transit and ensuring that only approved, privacy-compliant endpoints are used. Access policies and network segmentation must ensure that sensitive information does not leave the organizational perimeter unintentionally.

On-premises models offer the highest degree of control and visibility because all data, hardware, and operations remain entirely within the organization's private network perimeter. This approach comes with additional operational burdens, however, such as patching, hardware isolation, and local key management. Hybrid models—combining cloud inference with local orchestration—must account for secure communication and policy propagation between environments.

Hosting options where customers load their own models into a managed cloud environment add complexity to the security perimeter. Enterprises deploying fine-tuned or proprietary models must enforce strict access policies to prevent model misuse or theft of the proprietary model weights and ensure that inference traffic adheres to organizational security controls. An LLM Mesh can provide a standard and secure way of integrating these models by enforcing policy boundaries, propagating authorization metadata, and ensuring that fine-tuned models are only hosted and invoked in approved computing environments. This centralized enforcement significantly reduces the risk of accidental exposure or shadow deployments.

An LLM Mesh provides a unified architecture for enforcing consistent security policies across these diverse hosting scenarios, enabling teams to select deployment models based on business needs without sacrificing control or visibility.

## Isolation in Multitenant and Multiagent Environments

In agentic applications built on an LLM Mesh, isolation is critical to prevent cross-tenant data leakage, resource contention, or unauthorized lateral movement between agents or services. Because many services and objects may be shared across teams or tenants, enforcing both network-level and compute-level boundaries is essential.

*Network-level isolation* is typically achieved through mechanisms such as virtual private clouds (VPCs), dedicated subnets, and firewall rules that segment traffic between services. These controls prevent unauthorized communication between agents or tools that reside in different logical or organizational boundaries.

*Resource isolation* ensures that even when services share infrastructure, their execution environments remain logically separated. This is enforced using container technologies such as Docker, orchestration platforms like Kubernetes, and in high-assurance environments, hardware-backed secure enclaves (e.g., Intel SGX), which use encrypted memory regions within a CPU to protect code and data even from the host system. These techniques mitigate the risk of one agent affecting the behavior or visibility of another and are especially important in multiagent workflows that operate concurrently or on shared infrastructure.

An LLM Mesh enables consistent enforcement of these isolation strategies by standardizing how objects are provisioned, connected, and governed—regardless of how or where they are deployed.

## Multiagent Applications with Security Boundaries

Multiagent applications—where multiple agents collaborate to complete a task by exchanging outputs and invoking shared tools—are a powerful design pattern, but they introduce additional security risks. Without proper safeguards, one agent's output can become a vector for injecting prompts, introducing malformed content, or misclassifying sensitive data that is then consumed without verification by another agent or tool. These gaps can lead to unintended actions, unauthorized access, or compliance violations.

To mitigate these risks, every data exchange in a multiagent system must be treated as a potential attack surface. A key architectural principle is output validation: before content produced by one agent is handed off to another component, it must be checked for

structural integrity, adherence to policy boundaries, and metadata consistency. This includes scanning for prompt injection attempts, ensuring sensitive data classifications are respected, and verifying that content conforms to expected formats and schemas.

As introduced in [Chapter 5](#), moderation and safety filtering are essential safeguards for managing unsafe or hallucinated model outputs. Here, we extend that principle beyond safety into architectural enforcement: validation must occur at every handoff, not just at the application edge. An LLM Mesh enables this by embedding policy-aware validation layers throughout the system, ensuring that multiagent interactions remain secure, interpretable, and compliant.

## Aligning Agentic Application Security with Enterprise Standards

For enterprise adoption, an LLM Mesh must ensure that the security capabilities of the agentic applications built within it align with widely recognized governance frameworks and regulatory standards. While most of these standards—such as the NIST AI Risk Management Framework (NIST AI RMF), ISO/IEC 42001, and the European Union Agency for Cybersecurity (ENISA) guidelines—are voluntary, they offer a clear, structured path for organizations seeking to operationalize AI responsibly. An LLM Mesh supports this goal by making it possible to implement and enforce technical controls that map directly to these frameworks. For organizations aiming to meet these benchmarks, an LLM Mesh architecture becomes a foundational enabler—allowing high-level governance principles to be translated into consistent, enforceable practices.

Configuration management systems—such as infrastructure-as-code platforms or orchestration tools—allow teams to define policies in one place and deploy them uniformly across environments. These policies might govern access control, logging behavior, secrets handling, or audit retention.

As established in [Chapter 2](#), the LLM Mesh catalog plays a pivotal role in maintaining object metadata, version history, and lifecycle state. In this chapter, we see that the same catalog serves as a foundation for enforcing governance and security policies: ensuring that only authorized agents and tools are used, validating object compatibility, and enabling consistent enforcement of lifecycle-based

controls (e.g., deprecating objects that fail compliance checks with graceful fallbacks and alerting mechanisms).

This combination of declarative policy definition and catalog-backed enforcement is what allows an LLM Mesh to operationalize governance—not merely as documentation, but as a live, enforceable part of the infrastructure. By integrating auditability, policy enforcement, and lifecycle oversight directly into the architecture, an LLM Mesh allows the scalable development and deployment of agentic applications that meet enterprise security requirements.

## Conclusion

Unlike traditional systems where security measures are distributed across entities, an LLM Mesh ensures that enforcement occurs at every interaction boundary—between agents, tools, and services—with full awareness of context, identity, and metadata.

This architectural consistency enables organizations not only to implement best practices but also to prove that they are enforced. Structured logging, real-time anomaly detection, and output validation are embedded into the infrastructure, transforming observability from a diagnostic tool into a security and compliance guarantee.

Throughout this chapter, we’ve seen how an LLM Mesh supports fine-grained ABAC, enables federated identity across distributed systems, and aligns with leading governance frameworks like NIST AI RMF and ISO/IEC 42001. These capabilities allow enterprises to adopt AI at scale while staying in control of risk and compliance.

**Chapter 7** will conclude the book by illustrating how the principles and mechanisms covered in the previous chapters come together in a realistic enterprise environment—highlighting how an LLM Mesh enables secure, scalable, and governed AI deployment in practice.

# An LLM Mesh in Action: An Example Agentic Application

This concluding chapter serves as a capstone, moving from the *what* of an LLM Mesh to a complete, end-to-end *why*. The preceding chapters have established the principles, objects, and federated services that define an LLM Mesh. Now, we will build a multiagent application to demonstrate how an LLM Mesh enables you to create sophisticated and secure agents to solve a complex, high-stakes business problem.

Our case study is a multiagent Anti-Money Laundering (AML) Investigation Assistant. This system is designed not to replace, but to augment the capabilities of a human financial crime investigator, making their work faster, more thorough, and more effective. By examining its architecture, we will see how an LLM Mesh provides the necessary foundation for building powerful, reliable, and governed agentic applications in the enterprise.

We will first define the business problem and the high-pressure environment of AML compliance. Then, we will walk through the application's design and operation through the specific lens of each preceding chapter, connecting abstract concepts to concrete implementation choices.

# The Business Problem: Taming Complexity in Financial Crime Detection

Financial institutions are legally required to monitor transactions for suspicious activity and report their findings to government bodies to help combat financial crime. This critical responsibility falls to the AML Investigator, a role that combines forensic accounting with detective work. The core challenge this position faces is an overwhelming deluge of data. A single system-generated alert can require an investigator to manually sift through disparate sources: internal transaction records, customer identity documents, private case notes, and public news archives. The process is slow, laborious, and fraught with the risk of missing a critical connection. Given that fines for noncompliance can run into the billions, the stakes could not be higher.

To understand the business-specific vocabulary of the solution, a few key terms are essential:

## *Anti-Money Laundering (AML)*

The set of laws, regulations, and procedures designed to prevent criminals from disguising illegally obtained funds as legitimate income.

## *Know Your Customer (KYC)*

The mandatory process of verifying a client's identity and background, with the resulting data stored in internal records.

## *Politically exposed person (PEP)*

An individual holding a prominent public function (e.g., a government official) who is considered to pose a higher risk for potential involvement in bribery and corruption.

## *Adverse media*

Negative information discovered about a customer from public sources, such as news articles detailing involvement in crime or scandal.

## *Suspicious activity report (SAR)*

The official report an investigator must file with a financial crime authority, such as the Financial Crimes Enforcement Network (FinCEN) in the United States, if they suspect a client is engaged in illicit activities.

The solution presented here is an AML Investigation Assistant, which will be accessed through an AML Investigation platform. The AML Investigation platform is a web-based interface that provides access to authenticated users to the AML Investigation Assistant and other tools. It is not part of the AML Investigation Assistant, but it is the means by which it is accessed.

The AML Investigation Assistant acts as a copilot for the human investigator. When a new alert is triggered, the system automatically dispatches a team of specialized agents to conduct the initial research. There are five agents that will be part of our ML Investigation Assistant:

- The Alert Data Search Agent will gather all internal KYC data and past transaction history for the customer in question.
- The PEP Discovery Agent will scan external watchlists and public news for adverse media related to the customer.
- The Graph Exploration Agent will analyze the customer's network of transactions to visualize and identify hidden or unusual connections to other entities.
- The AML Assistant Agent, acting as an orchestrator, will synthesize the findings from the other agents into a single, preliminary investigation report.
- If the initial findings warrant it, a SAR Preparation Agent will use the synthesized report to draft a regulatory-compliant SAR.

The human investigator receives this comprehensive report within minutes of the initial alert. This frees them from hours of manual data collection, allowing them to apply their expertise to what matters most: exercising judgment, making the final decision, and protecting the institution from risk.

**Figure 7-1** illustrates the structure of this multiagent application.

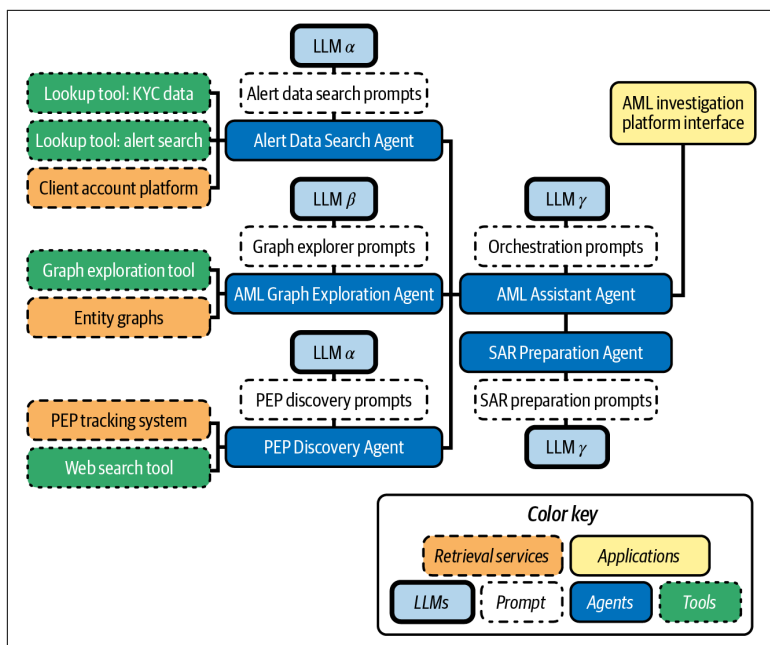


Figure 7-1. Structure of AML Investigation Assistant using multiple, different LLMs for different purposes

## Building Within an LLM Mesh Architecture: A Practical Walkthrough

The following subsections describe how this multiagent application would be built in an LLM Mesh, with each section referencing a previous chapter in the book.

### The Right Model for the Right Application (Chapter 1)

A one-size-fits-all model strategy would be both inefficient and ineffective for the AML Assistant. An LLM Mesh architecture allows the application to leverage different models for different tasks, striking an optimal balance between performance and cost. Here are the two strategies that guide the development of our app:

#### *A mixed-model strategy*

The application employs a mix of powerful and specialized models. The central AML Assistant Agent, which must synthesize complex, unstructured findings, and the SAR Preparation Agent,

which must generate coherent, legally sound narratives, both use a powerful, state-of-the-art model (in our case, GPT-4o). The accuracy and reasoning required for these tasks justify the higher cost. In contrast, the Alert Data Search Agent, which performs narrower tasks like extracting specific entities from KYC files, uses a smaller, faster, and potentially cheaper model (in our case, Llama 3.1 8B).

#### *A hybrid hosting strategy*

An LLM Mesh enables a hybrid hosting model to meet security and cost requirements. The models processing sensitive internal data (transaction records, KYC files) are run on a self-hosted LLM service inside a secure VPC. This ensures no customer data leaves the organizational firewall. Conversely, the PEP Discovery Agent, which only interacts with public information, uses a third-party MaaS provider and is given a tool to call a web search API. This approach optimizes for cost without compromising data security.

## **The Objects of the AML Assistant (Chapter 2)**

The entire application is composed of modular, reusable objects, each registered and documented in the central catalog of an LLM Mesh. This composability is what makes the system manageable. Here are the objects of our assistant:

#### *LLMs and LLM services*

The agents are powered by the specific LLMs (e.g., GPT-o3, Llama-3.1-8b-instruct), which are accessed via the API endpoints for the self-hosted and third-party LLM services. These models and services are registered in the catalog, and permissions for the use of the different models and services are recorded. For example, the third-party LLM services are not allowed for the most sensitive data.

#### *Agents*

The five distinct agents (AML Assistant, Alert Data Search, etc.) are cataloged objects. Their entries include descriptions of their capabilities, version history, and dependencies.

#### *Tools*

The agents use a variety of cataloged tools, each with a defined schema explaining its function and expected inputs. These

include an internal transaction database (a data querying service), a PEP Watchlist API and a web search tool (both API services), and a graph visualization tool.

#### *Prompts*

Prompts are treated as version-controlled assets. The SAR Preparation Agent uses a highly structured prompt template with placeholders for “Subject Name” and “Suspicious Activity Description,” ensuring its output format is always compliant with regulatory filing standards.

#### *Retrieval services*

The application includes a retrieval service over a vector store containing the bank’s internal AML policies and thousands of past case files. This allows the AML Assistant Agent to ask, “Does this activity pattern match known typologies of money laundering?”

#### *Application*

The AML Investigation platform is the user-facing application object. It provides the UI where the human investigator reviews the generated report, converses with the AML Assistant Agent for follow-up questions, and finalizes the SAR for submission.

## **Optimizing Cost in a High-Stakes Environment** **(Chapter 3)**

In a function as cost-sensitive as compliance, managing spend is critical. These federated services of an LLM Mesh provide the necessary visibility and control:

#### *Cost visibility*

The centralized cost service tracks spend per investigation, per agent, and per tool. The compliance department gets a clear, auditable report on the operational cost of its AML program, allowing it to justify budgets and identify inefficiencies.

#### *Cost reduction techniques*

Several automated techniques are applied. The high-volume queries made by the Alert Data Search Agent are a prime candidate for model substitution. By using a smaller, self-hosted model that meets the quality bar, the organization takes on a significant fixed hardware cost. However, for a task with high, predictable volume, this approach is more cost-effective

than using a larger model or a pay-per-use service. When the agents retrieve long documents to be summarized, prompt compression is automatically applied to reduce the input token count. Finally, for the repetitive system prompts that define each agent's role, the system leverages the LLM service's context caching to reduce costs on every run.

## Measuring Performance and Quality (Chapter 4)

For the AML Assistant, performance is not just about speed; it's about accuracy and reliability. The key business metric is the reduction in *time to decision* for the human investigator, without any degradation in the quality of those decisions. An LLM Mesh supports both intrinsic and extrinsic quality monitoring to measure that we've met these goals:

### *Intrinsic quality monitoring*

Before deployment, the consistency of each agent is tested. The SAR Preparation Agent, for example, is run hundreds of times on a golden dataset of test cases to ensure its output structure remains stable and predictable. The retrieval quality of the internal knowledge base is also continually monitored to ensure it surfaces the most relevant policy documents.

### *Extrinsic quality monitoring*

The system relies on two layers of extrinsic evaluation. First, an automated LLM-as-a-judge agent checks every drafted SAR against the source documents for factuality. Any claim in the draft that cannot be directly attributed to a source is flagged for human review. Second, human feedback provides the ultimate quality signal. The investigator's final edits to the SAR draft are captured and used as high-quality training data to continually fine-tune the SAR Preparation Agent, aligning its output ever more closely with expert expectations.

## Ensuring Safety and Reliability (Chapter 5)

In a regulatory context, an unsafe output—such as a hallucinated or biased statement—is not just an error; it's a critical compliance failure. Building an agentic application in an LLM Mesh enables developers to embed safety throughout the assistant's architecture. Here are three critical safety factors that our LLM Mesh will provide to ensure the safety of the application:

### *Hallucination prevention*

An unsubstantiated claim in a SAR could lead to legal action. To reduce this risk, the application is built on a strict “grounding” principle. Every piece of information in the final report is designed to be explicitly sourced and hyperlinked back to the internal document or external web page it came from, allowing for human verification. The agents are instructed not to infer information but only to report on and synthesize verifiable facts from the provided sources.

### *Filtering and adversarial testing*

A PII filter is automatically applied before any query is sent to the external web search tool, preventing accidental leakage of customer data. A forbidden content filter is in place to prevent the agents from providing legal or investment advice, which they are not qualified to do. An adversarial testing suite runs nightly, using a library of prompts designed to trick the agents into revealing confidential information or bypassing controls, allowing vulnerabilities to be patched proactively.

### *Human-in-the-loop (HITL) workflow*

The UI is designed for explainability. The investigator can click on any statement in the summary and see the full chain of evidence—which agent found it, from what source, and the exact text. The entire system is designed as a HITL workflow; the agent drafts, but the human decides and remains fully accountable.

## **Securing a Multiagent System (Chapter 6)**

The AML Assistant interacts with both sensitive internal data and the public internet, making a robust security posture nonnegotiable. An LLM Mesh architecture should make it simple for builders of agentic applications to enforce security at every interaction boundary. Here are the most important ways that the LLM Mesh ensures security in our multiagent assistant:

### *Fine-grained access control*

Using ABAC, the system enforces the principle of least privilege at the object level. The Alert Data Search Agent is granted read-only access to specific internal databases but is explicitly denied access to any external network tools. Conversely, the PEP Discovery Agent can access the public web but is denied access to all internal customer databases.

### *Secure gateway and audit trail*

All interagent communication and tool calls are routed through the gateway of an LLM Mesh architecture. This creates a single, immutable, and end-to-end audit log for every action taken during an investigation. This verifiable trail is crucial for satisfying regulatory reviews.

### *Secrets management and isolation*

The API key for the external web search tool is managed by a centralized secrets manager and rotated automatically; it is never hardcoded. Furthermore, the agent responsible for handling internal PII runs in a separate, more secure compute environment than the agent that interacts with the public internet, preventing any potential for lateral movement or cross-contamination.

## **The Benefits of an LLM Mesh for Building, Deploying, and Maintaining the AML Assistant**

Having dissected the AML Assistant, it becomes clear that building it as a traditional, monolithic application would introduce substantial technical debt and operational risk. With an LLM Mesh architecture, builders of agentic applications can benefit from a series of compounding benefits that make such a system not just possible, but practical:

### *Speed—accelerating development and innovation*

By providing a standardized framework, the LLM Mesh has allowed the organization to move from ad hoc experimentation to a strategic, scalable capability. Without it, each development team would need to implement its own solutions for API calls, security checks, and logging actions. This would lead to duplicated effort, inconsistent security, and fragmented logs. With the LLM Mesh, these are centrally managed, freeing developers to focus on business logic. The result is faster development, lower maintenance costs, and a more robust, secure, and governable AI ecosystem.

### *Reliability—building a robust, enterprise-grade system*

A monolithic AML application would be brittle. If the third-party web search API changes, the entire application could break. In an LLM Mesh architecture, only the tool's adapter

needs to be updated, with no impact on the agents that use it. Furthermore, federated services for performance monitoring and safety are applied consistently. This means every agent and every tool is subject to the same rigorous standards for hallucination detection and quality evaluation, creating a system that is far more reliable and predictable than one assembled from inconsistently governed parts.

#### *Trust and governance—the foundation for compliance*

For an AML application, trust is not optional; its core components, like auditability and transparency, are regulatory requirements. An auditor must be able to verify the entire decision-making process. A monolithic application with fragmented logs and hardcoded permissions would fail this test. An LLM Mesh architecture, by routing every action through a secure gateway, allows for a centralized, immutable audit trail for every investigation. The fine-grained access controls and centrally enforced safety policies provide provable governance. This architectural approach is what makes the system trustworthy to internal stakeholders, compliance officers, and external regulators.

## **Conclusion: From Ad Hoc AI to a Strategic Capability**

An LLM Mesh architecture provides the necessary foundation for making agentic applications feasible. Its core tenets—abstraction through a unified gateway, governance through a central catalog, and consistency through federated services for cost, performance, safety, and security—are what allow an enterprise to overcome the complexity of building and managing agentic applications at scale. An LLM Mesh architecture enables an enterprise to transition from isolated, experimental AI projects to a scalable, efficient, and strategic “factory” for building and deploying a portfolio of agentic applications that drive real and defensible business value.

## About the Author

---

As Head of AI Strategy at Dataiku, **Kurt Muehmel** brings Dataiku's vision of AI Success to industry analysts and media worldwide. He advises Dataiku's C-Suite on market and technology trends, ensuring that they maintain their position as pioneers.

Kurt is a creative and analytical executive with 15+ years of experience and foundational expertise in the Enterprise AI space and, more broadly, B2B SaaS go-to-market strategy and tactics. He's focused on building a future where the most powerful technologies serve the needs of people and businesses.